

1 **Predicting wildlife hosts of betacoronaviruses for SARS-CoV-2 sampling prioritization**

2

3 Daniel J. Becker<sup>1,†</sup>, Gregory F. Albery<sup>2,†</sup>, Anna R. Sjodin<sup>3</sup>, Timothée Poisot<sup>4</sup>, Tad A. Dallas<sup>5</sup>, Evan  
4 A. Eskew<sup>6,7</sup>, Maxwell J. Farrell<sup>8</sup>, Sarah Guth<sup>9</sup>, Barbara A. Han<sup>10</sup>, Nancy B. Simmons<sup>11</sup>, and Colin J.  
5 Carlson<sup>12,13,\*</sup>

6

7

8

9 † These authors share lead author status

10 \* Corresponding author: [colin.carlson@georgetown.edu](mailto:colin.carlson@georgetown.edu)

11

12 1. Department of Biology, Indiana University, Bloomington, IN, U.S.A.

13 2. Department of Biology, Georgetown University, Washington, D.C., U.S.A.

14 3. Department of Biological Sciences, University of Idaho, Moscow, ID, U.S.A.

15 4. Université de Montréal, Département de Sciences Biologiques, Montréal, QC, Canada.

16 5. Department of Biological Sciences, Louisiana State University, Baton Rouge, LA, U.S.A.

17 6. Department of Ecology, Evolution, and Natural Resources, Rutgers University, New Brunswick,  
18 NJ, U.S.A.

19 7. Department of Biology, Pacific Lutheran University, Tacoma, WA, U.S.A.

20 8. Department of Ecology & Evolutionary Biology, University of Toronto, Toronto, ON, Canada.

21 9. Department of Integrative Biology, University of California Berkeley, Berkeley, CA, U.S.A.

22 10. Cary Institute of Ecosystem Studies, Millbrook, NY, U.S.A.

23 11. Department of Mammalogy, Division of Vertebrate Zoology, American Museum of Natural  
24 History, New York, NY, U.S.A.

25 12. Center for Global Health Science and Security, Georgetown University Medical Center,  
26 Washington, D.C., U.S.A.

27 13. Department of Microbiology and Immunology, Georgetown University Medical Center,  
28 Washington, D.C., U.S.A.

29 **Abstract.**

30

31 Despite massive investment in research on reservoirs of emerging pathogens, it remains  
32 difficult to rapidly identify the wildlife origins of novel zoonotic viruses. Viral surveillance  
33 is costly but rarely optimized using model-guided prioritization strategies, and predictions  
34 from a single model may be highly uncertain. Here, we generate an ensemble of seven  
35 network- and trait-based statistical models that predict mammal-virus associations, and  
36 we use model predictions to develop a set of priority recommendations for sampling  
37 potential bat reservoirs and intermediate hosts for SARS-CoV-2 and related  
38 betacoronaviruses. We find nearly 300 bat species globally could be undetected hosts of  
39 betacoronaviruses. Although over a dozen species of Asian horseshoe bats (*Rhinolophus*  
40 spp.) are known to harbor SARS-like viruses, we find at least two thirds of betacoronavirus  
41 reservoirs in this bat genus might still be undetected. Although identification of other  
42 probable mammal reservoirs is likely beyond existing predictive capacity, some of our  
43 findings are surprisingly plausible; for example, several civet and pangolin species were  
44 highlighted as high-priority species for viral sampling. Our results should not be over-  
45 interpreted as novel information about the plausibility or likelihood of SARS-CoV-2's  
46 ultimate origin, but rather these predictions could help guide sampling for novel  
47 potentially zoonotic viruses; immunological research to characterize key receptors (e.g.,  
48 ACE2) and identify mechanisms of viral tolerance; and experimental infections to quantify  
49 competence of suspected host species.

## Main text.

50  
51  
52 Coronaviruses are a diverse family of positive-sense, single-stranded RNA viruses, found widely  
53 in mammals and birds<sup>1</sup>. They have a broad host range, a high mutation rate, and the largest  
54 genomes of any RNA viruses, but they have also evolved mechanisms for RNA proofreading and  
55 repair, which help to mitigate the deleterious effects of a high recombination rate acting over a  
56 large genome<sup>2</sup>. Consequently, coronaviruses fit the profile of viruses with high zoonotic potential.  
57 There are seven human coronaviruses (two in the genus *Alphacoronavirus* and five in  
58 *Betacoronavirus*), of which three are highly pathogenic in humans: SARS-CoV, SARS-CoV-2, and  
59 MERS-CoV. These three are zoonotic and widely agreed to have evolutionary origins in bats<sup>3-6</sup>.  
60  
61 Our collective experience with both SARS-CoV and MERS-CoV illustrate the difficulty of tracing  
62 specific animal hosts of emerging coronaviruses. During the 2002–2003 SARS epidemic, SARS-  
63 CoV was traced to the masked palm civet (*Paguma larvata*)<sup>7</sup>, but the ultimate origin remained  
64 unknown for several years. Horseshoe bats (family Rhinolophidae: *Rhinolophus*) were implicated  
65 as reservoir hosts in 2005, but their SARS-like viruses were not identical to circulating human  
66 strains<sup>4</sup>. Stronger evidence from 2017 placed the most likely evolutionary origin of SARS-CoV in  
67 *Rhinolophus ferrumequinum* or potentially *R. sinicus*<sup>8</sup>. Presently, there is even less certainty in the  
68 origins of MERS-CoV, although spillover to humans occurs relatively often through contact with  
69 dromedary camels (*Camelus dromedarius*). A virus with 100% nucleotide identity in a ~200 base  
70 pair region of the polymerase gene was detected in *Taphozous* bats (family Emballonuridae) in  
71 Saudi Arabia<sup>9</sup>; however, based on spike gene similarity, other sources treat HKU4 virus from  
72 *Tylonycteris* bats (family Vespertilionidae) in China as the closest-related bat virus<sup>10,11</sup>. Several  
73 bat coronaviruses have shown close relation to MERS-CoV, with a surprisingly broad geographic  
74 distribution from Mexico to China<sup>12,13,14,15</sup>.  
75  
76 Coronavirus disease 2019 (COVID-19) is caused by severe acute respiratory syndrome  
77 coronavirus-2 (SARS-CoV-2), a novel virus with presumed evolutionary origins in bats. Although  
78 the earliest cases were linked to a wildlife market, contact tracing was limited, and there has been  
79 no definitive identification of the wildlife contact that resulted in spillover nor a true “index case.”  
80 Two bat viruses are closely related to SARS-CoV-2: RaTG13 bat CoV from *Rhinolophus affinis*  
81 (96% identical overall), and RmYN02 bat CoV from *Rhinolophus malayanus* (97% identical in one  
82 gene but only 61% in the receptor binding domain and with less overall similarity)<sup>6,16</sup>. The  
83 divergence time between these bat viruses and human SARS-CoV-2 has been estimated as 30-70  
84 years<sup>17</sup>, suggesting that the main host(s) involved in spillover remain unknown. Evidence of viral  
85 recombination in pangolins has been proposed but is unresolved<sup>17</sup>. SARS-like betacoronaviruses  
86 have been recently isolated from Malayan pangolins (*Manis javanica*) traded in wildlife  
87 markets<sup>18,19</sup>, and these viruses have a very high amino acid identity to SARS-CoV-2, but only show  
88 a ~90% nucleotide identity with SARS-CoV-2 or Bat-CoV RaTG13<sup>20</sup>. None of these host species  
89 are universally accepted as the origin of SARS-CoV-2 or a progenitor virus, and a “better fit” wildlife  
90 reservoir could likely still be identified. However, substantial gaps in betacoronavirus sampling

91 across wildlife limit actionable inference about plausible reservoirs and intermediate hosts for  
92 SARS-CoV-2<sup>21</sup>.

93  
94 Identifying likely reservoirs of zoonotic pathogens is challenging<sup>22</sup>. Sampling wildlife for the  
95 presence of active or previous infection (i.e., seropositivity) represents the first stage of a pipeline  
96 for proper inference of host species<sup>23</sup>, but sampling is often limited in phylogenetic, temporal, and  
97 spatial scale by logistical constraints<sup>24</sup>. Given such restrictions, modeling efforts can play a  
98 critical role in helping to prioritize pathogen surveillance by narrowing the set of plausible  
99 sampling targets<sup>25</sup>. For example, machine learning approaches have generated candidate lists of  
100 likely, but unsampled, primate reservoirs for Zika virus, bat reservoirs for filoviruses, and avian  
101 reservoirs for *Borrelia burgdorferi*<sup>26–28</sup>. In some contexts, models may be more useful for  
102 identifying which host or pathogen groups are *unlikely* to have zoonotic potential<sup>29</sup>. However,  
103 these approaches are generally applied individually to generate predictions. Implementation of  
104 multiple modeling approaches collaboratively and simultaneously could reduce redundancy and  
105 apparent disagreement at the earliest stages of pathogen tracing and help advance modeling  
106 work by addressing inter-model reliability, predictive accuracy, and the broader utility (or  
107 inefficacy) of such models in zoonosis research.

108  
109 Because SARS-like viruses (subgenus *Sarbecovirus*) are only characterized from a small number  
110 of bat species in publicly available data, current modeling methods are poorly tailored to exactly  
111 infer their potential reservoir hosts. In this study, we instead conduct two predictive efforts that  
112 may help guide the inevitable search for known and future zoonotic coronaviruses in wildlife: (1)  
113 broadly identifying bats and other mammals that may host any *Betacoronavirus* and (2)  
114 specifically identifying species with a high viral sharing probability with the two *Rhinolophus*  
115 species carrying the closest known wildlife relatives of SARS-CoV-2. To do this, we developed a  
116 standardized dataset of mammal-virus associations by integrating a previously published  
117 mammal-virus dataset<sup>30</sup> with a targeted scrape of all GenBank coronavirus accessions and their  
118 associated hosts. Our final dataset spanned 710 host species and 359 virus genera, including  
119 107 mammal hosts of betacoronaviruses as well as hundreds of other (non-coronavirus)  
120 association records. We harmonized our host-virus data with a mammal phylogenetic supertree<sup>31</sup>  
121 and over 60 ecological traits of bat species<sup>27,32,33</sup>. Using these standardized data, six subteams  
122 generated seven predictive models of host-virus associations, including four network-based and  
123 three trait-based approaches. These efforts generated seven ranked lists of suspected bat hosts  
124 of betacoronaviruses and five ranked lists for other mammals. Each ranked list was scaled  
125 proportionally and consolidated in an ensemble of recommendations for betacoronavirus  
126 sampling and broader eco-evolutionary research (ED Figure 1).

127  
128 In our ensemble, we draw on two popular approaches to identify candidate reservoirs and  
129 intermediate hosts of betacoronaviruses. *Network-based methods* estimate a full set of “true”  
130 unobserved host-virus interactions based on a recorded network of associations (here, pairs of  
131 host species and associated viral genera). These methods are increasingly popular as a way to  
132 identify latent processes structuring ecological networks<sup>34–36</sup>, but they are often confounded by

133 sampling bias and can only make predictions for species within the observed network (i.e., those  
134 that have available virus data; in-sample prediction). In contrast, *trait-based methods* use  
135 observed relationships concerning host traits to identify species that fit the morphological,  
136 ecological, and/or phylogenetic profile of known host species of a given pathogen and rank the  
137 suitability of unknown hosts based on these trait profiles<sup>28,37</sup>. These methods may be more likely  
138 to recapitulate patterns in observed host-pathogen association data (e.g., geographic biases in  
139 sampling, phylogenetic similarity in host morphology), but they more easily correct for sampling  
140 bias and can predict host species without known viral associations (out-of-sample prediction).

141  
142 Predictions of bat betacoronavirus hosts derived from network- and trait-based approaches  
143 displayed strong inter-model agreement within-group, but less with each other (Figure 1A,B). In-  
144 sample, we identified bat species across a range of genera as having the highest predicted  
145 probabilities of hosting betacoronaviruses, distributed in distinct families in both the Old World  
146 (e.g., Hipposideridae, several subfamilies in the Vespertilionidae) and the New World (e.g.,  
147 *Artibeus jamaicensis* from the Phyllostomidae; Figure 1C). Out-of-sample, our multi-model  
148 ensemble more conservatively limited predictions to primarily Old World families such as  
149 Rhinolophidae and Pteropodidae (Figure 1D). Of the 1,037 bat hosts not currently known to host  
150 betacoronaviruses, our models identified between 1 and 720 potential hosts based on a 10%  
151 omission threshold (90% sensitivity). Applying this same threshold to our ensemble predictions,  
152 we identified 291 bat species that are likely undetected hosts of betacoronaviruses. These  
153 include approximately half of bat species in the genus *Rhinolophus* not currently known to be  
154 betacoronavirus hosts (30 of 61), compared to 16 known hosts in this genus. Given known roles  
155 of rhinolophids as hosts of SARS-like viruses, our results suggest that SARS-like virus diversity  
156 could be undescribed for around two-thirds of the potential reservoir bat species.

157  
158 Our multi-model ensemble predicted undiscovered betacoronavirus bat hosts with striking  
159 geographic patterning (Figure 2). In-sample, the top 50 predicted bat hosts were broadly  
160 distributed and recapitulated observed patterns of bat betacoronavirus hosts in Europe, parts of  
161 sub-Saharan Africa, and southeast Asia, although our models also predicted greater-than-  
162 expected richness of likely bat reservoirs in the Neotropics and North America. In contrast, the  
163 top out-of-sample predictions clustered in Vietnam, Myanmar, and southern China.

164  
165 Because only trait-based models were capable of out-of-sample prediction, the differences in  
166 geographic patterns of our predictions likely reflect distinctions between the network- and trait-  
167 based modeling approaches, which we suggest should be considered qualitatively different lines  
168 of evidence. Network approaches proportionally upweight species with high observed viral  
169 diversity, recapitulating sampling biases largely unrelated to coronaviruses (e.g., frequent  
170 screening for rabies lyssaviruses in vampire bats, which have been sampled in a comparatively  
171 limited capacity for coronaviruses<sup>14,38-40</sup>). Highly ranked species may also have been previously  
172 sampled without evidence of betacoronavirus presence; for example, *Rhinolophus luctus* and  
173 *Macroglossus sobrinus* from China and Thailand, respectively, tested negative for  
174 betacoronaviruses, but detection probability was limited by small sample sizes<sup>41-43</sup>. In contrast,

175 trait-based approaches are constrained by their reliance on phylogeny and ecological traits, and  
176 the use of geographic covariates made models more likely to recapitulate existing spatial  
177 patterns of betacoronavirus detection (i.e., clustering in southeast Asia). However, their out-of-  
178 sample predictions are, by definition, inclusive of unsampled hosts<sup>44</sup>, which potentially offer  
179 greater return on viral discovery investment.

180  
181 Multi-model ensemble predictions also clustered taxonomically along parallel lines. Applying a  
182 graph partitioning algorithm (phylogenetic factorization) to the bat phylogeny<sup>45</sup>, we found that in-  
183 sample predictions were on average lowest for the Yangochiroptera (Figure 3). This makes  
184 intuitive sense, because this clade does not include the groups known to harbor the majority of  
185 betacoronaviruses detected in bats (e.g., *Rhinolophus*, Hipposideridae). Out-of-sample  
186 predictions were lower in the New World superfamily Noctilionoidea and the emballonurids,  
187 whereas several subfamilies of Old World fruit bats<sup>46</sup>, including the Rousettinae, Cynopterinae,  
188 and Eidolinaei, had higher mean probabilities of betacoronavirus hosting. Lastly, our ensemble  
189 also identified the *Rhinolophus* genus as having greater mean probabilities (ED Table 1).

190  
191 These clade-specific patterns of predicted probabilities across extant bats could be particularly  
192 applicable for guiding future surveillance. On the one hand, betacoronavirus sampling in  
193 southeast Asian bat taxa (especially the genus *Rhinolophus*) may have a high success of viral  
194 detection but may not improve existing bat sampling gaps<sup>47</sup>. On the other hand, discovery of novel  
195 betacoronaviruses in Neotropical bats or Old World fruit bats could significantly revise our  
196 understanding of the bat-virus association network. Such discoveries would be particularly  
197 important for global health security, given the surprising identification of a MERS-like virus in  
198 Mexican bats<sup>14</sup> and the likelihood that post-COVID pandemic preparedness efforts will focus  
199 disproportionately on Asia despite the near-global presence of bat betacoronaviruses.

200  
201 Although our ensemble model of potential bat betacoronavirus reservoirs generated strong and  
202 actionable predictions, our mammal-wide predictions were largely uninformative. In particular,  
203 minimal inter-model agreement (ED Figure 2) indicated a lack of consistent, biologically  
204 meaningful findings. Major effects of sampling bias were apparent from the top-ranked species,  
205 which were primarily domestic animals or well-studied mesocarnivores (ED Figure 2B).  
206 Phylogenetic factorization mostly failed to find specific patterns in prediction (ED Table 2): in-  
207 sample, mean predictions primarily confirmed betacoronavirus detection in the remaining  
208 Laurasiatheria (e.g., ungulates, carnivores, pangolins, hedgehogs, shrews), although nested  
209 clades of marine mammals (i.e., cetaceans) were less likely to harbor these viruses, as expected  
210 given betacoronavirus epidemiology and their predominance in terrestrial mammals. Our  
211 mammal predictions thus reflect a combination of detection bias and poor performance of  
212 network methods on limited data that likely signals the limits of existing predictive capacity. Our  
213 dataset contained only 30 non-bat betacoronavirus hosts, many of which were identified during  
214 sampling efforts following the first SARS outbreak<sup>7</sup>. Although the laurasiatherians are likely to  
215 include more potential intermediate hosts than other mammals, the high diversity of this clade  
216 restricts insights for sampling prioritization, experimental work, or spillover risk management.



217

218 Given the unresolved origins of SARS-CoV-2 and significant motivation to identify other SARS-like  
219 coronaviruses and their reservoir hosts for pandemic preparedness<sup>21</sup>, we further explored our  
220 only model that could generate out-of-sample predictions for all mammals<sup>48</sup>. This model uses  
221 geographic distributions and phylogenetic relatedness to estimate viral sharing probability.  
222 Where one or more (potential) hosts are known, these sharing patterns can be interpreted to  
223 identify probable reservoir hosts<sup>48</sup>. Because *Rhinolophus affinis* and *R. malayanus* host viruses  
224 that are closely related to SARS-CoV-2<sup>6,16</sup>, we used their predicted sharing patterns to identify  
225 possible reservoirs of sarbecoviruses. In doing so, we aimed to work around a major data  
226 limitation: fewer than 20 sarbecovirus hosts were recorded in our dataset, a sample size that  
227 would preclude most modeling approaches.

228

229 For both presumed bat host species of sarbecoviruses, the most probable viral sharing hosts  
230 were again within the Laurasiatheria. Although bats—especially rhinolophids—unsurprisingly  
231 assumed the top predictions given phylogenetic affinity with known hosts (ED Table 3, ED Figure  
232 3), several notable patterns emerged in the rankings of other mammals. Pangolins (Pholidota)  
233 were disproportionately likely to share viruses with *R. affinis* and *R. malayanus* (ED Figure 4); the  
234 Sunda pangolin (*Manis javanica*) and Chinese pangolin (*M. pentadactyla*) were in the top 20  
235 predictions for both reservoir species (ED Table 4). This result is promising given the much-  
236 discussed discovery of SARS-like betacoronaviruses in *M. javanica*<sup>18</sup>. The Viverridae were also  
237 disproportionately well-represented in the top predictions (ED Figure 5), most notably the masked  
238 palm civet (*Paguma larvata*), which was identified as an intermediate host of SARS-CoV<sup>49,50</sup> (ED  
239 Table 4).

240

241 The ability of our virus sharing model to capture known patterns of coronavirus hosts using only  
242 two predictor variables is encouraging, and implies that mammal phylogeography has played a  
243 predictable role in historical betacoronavirus spillover. Moreover, these findings lend credibility  
244 to other predictions of SARS-CoV-2 sharing patterns and host susceptibility. Many of the model's  
245 top predictions were mustelids (i.e., ferrets and weasels), and the most likely viral sharing partner  
246 for both *Rhinolophus* species was the hog badger (*Arctonyx collaris*; ED Table 4). Taken together  
247 with reports of SARS-CoV-2 spread in mink farms<sup>51</sup>, these results highlight the relatively  
248 unexplored potential for mustelids to serve as betacoronavirus hosts. Similarly, identification of  
249 several deer and Old World monkey taxa as high-probability hosts in our clade-based analysis (ED  
250 Figure 3) meshes with the observation of high binding of SARS-CoV-2 to ACE2 receptors in cervid  
251 deer and primates<sup>52</sup>. Felids (especially leopards) also ranked relatively high in our viral sharing  
252 predictions (ED Table 4, ED Figure 5), which is of particular interest given reports of SARS-CoV-2  
253 susceptibility among cats<sup>53</sup>. However, we caution that this model was the only approach in our  
254 ensemble that could generate out-of-sample prediction across mammals, and therefore its  
255 predictions lacked confirmation (and filtering of potential spurious results) by other models that  
256 were designed and implemented independently.

257

258 Several limitations apply to our work, most notably the difficulty of empirically verifying  
259 predictions. Although some virological studies have incidentally tested specific hypotheses (e.g.,  
260 filovirus models and bat surveys<sup>27,54</sup>, henipavirus models and experimental infections<sup>23,55</sup>), model-  
261 based predictions are nearly never subject to systematic verification or post-hoc efforts to identify  
262 and correct spurious results. Greater dialogue between modelers and empiricists is necessary to  
263 systematically confront the growing set of predicted host-virus associations with experimental  
264 validation or field observation. *Scotophilus heathii*, *Hipposideros larvatus*, and *Pteropus lylei*, all  
265 highly predicted bat species in our out-of-sample rankings, have been reported positive for  
266 betacoronaviruses in the literature<sup>43,56</sup>; however, resulting sequences were not annotated to  
267 genus level in GenBank. These results support the idea that our models identified relevant targets  
268 correctly but also highlight an evident limitation of the workflow. Whereas an automated  
269 approach was the ideal method to systematically compile over 30,000 samples on the timescales  
270 commensurate with ongoing efforts to trace SARS-CoV-2 in wildlife, we suggest this discrepancy  
271 highlights the need for careful virological work downstream at every stage of the modeling  
272 process, including the development of hybrid manual-automated data pipelines.

273  
274 Additionally, overcoming underlying model biases that are driven by historical sampling regimes  
275 will require coordinated efforts in field study design. Bat sampling for betacoronaviruses has  
276 prioritized viral discovery<sup>39,40,57-59</sup>, but limitations in the spatial and temporal scale (and  
277 replication) of field sampling have likely created fundamental gaps in our understanding of  
278 infection dynamics in bat populations<sup>24</sup>. Limited longitudinal sampling of wild bats suggests  
279 betacoronavirus detection is sporadic over time and space<sup>56,60</sup>, implying strong seasonality in  
280 virus shedding pulses<sup>61</sup>. Carefully tailored spatial and temporal sampling efforts for priority taxa  
281 identified here, within the *Rhinolophus* genus or other high-prediction bat clades, will be key to  
282 identifying the environmental drivers of betacoronavirus shedding from wild bats and possible  
283 opportunities for contact between bats, intermediate hosts, and humans.

284  
285 Future field studies will undoubtedly be important to understand viral dynamics in bats but are  
286 inherently costly and labor-intensive. These efforts are particularly challenging during a pandemic,  
287 as many scientific operations have been suspended, including field studies of bats in some  
288 regions to limit possible viral spillback from humans. However, various alternative efforts could  
289 both advance basic virology and allow testing model predictions. General open access to viral  
290 association records, including GenBank accessions and the upcoming release of the USAID  
291 PREDICT program's data, could answer open questions and allow updates to our sampling  
292 prioritization (including potentially modeling at subgenus level, with greater data availability).  
293 Museum specimens and historical collections from diverse research programs also offer key  
294 opportunities to retrospectively screen samples from bats and other mammals for  
295 betacoronaviruses and to enhance our understanding of complex host-virus interactions<sup>62</sup>.  
296 Large-scale research networks, such as GBatNet (Global Union of Bat Diversity Networks) and its  
297 member networks, could provide diverse samples and ensure proper partnerships and equitable  
298 access and benefit sharing of knowledge across countries<sup>63,64</sup>. Whole-genome sequencing  
299 through initiatives such as the Bat1K Project (<https://bat1k.ucd.ie>) would facilitate fundamental



300 and applied insights into the immunological pathways through which bats can apparently harbor  
301 many virulent viruses (including but not limited to betacoronaviruses) without displaying clinical  
302 disease<sup>65,66</sup>.

303  
304 To expedite such work, we have made our binary predictions of host-virus associations for all  
305 seven models and all 1,000+ bat species publicly available (Supplementary Table 1). Such results  
306 are provided both in the spirit of open science and with the hope that future viral detection,  
307 isolation, or experimental studies might confirm some of these predictions or rule out others<sup>55</sup>. In  
308 ongoing collaborative efforts, we aim to consolidate results from field studies that address these  
309 predictions (e.g., serosurveys) and to track Genbank submissions to expand the known list of  
310 betacoronavirus hosts. In several years, we intend to revisit these predictions as a post-hoc test  
311 of model validation, which would represent the first effort to test the performance of such models  
312 and assess their contribution to basic science and to pandemic preparedness.

313  
314 It is crucial that our predictions be interpreted as a set of hypotheses about potential host-virus  
315 compatibility rather than strong evidence that a particular mammal species is a true reservoir for  
316 betacoronaviruses. In particular, susceptibility is only one aspect of host competence<sup>22,67</sup>, which  
317 encompasses the diverse genetic and immunological processes that mediate within-host  
318 responses following exposure<sup>68</sup>. SARS-CoV-2 in particular may have a broad host range<sup>52</sup>, given  
319 hypothesized compatibility with the ACE2 receptor in many mammal species, but this only adds  
320 to the extreme caution with which any data should be used to implicate a potential wildlife  
321 reservoir of the virus, given that rapid interpretation of inconclusive molecular evidence has likely  
322 already generated spurious reservoir identifications<sup>69,70</sup>. Future efforts to isolate live virus from  
323 wildlife or to experimentally show viral replication would more robustly test whether predicted  
324 host species actually play a role in betacoronavirus maintenance in wildlife<sup>55</sup>.

325  
326 Without direct lines of virological evidence, we note that our sampling prioritization scheme also  
327 does not implicate any given mammal species in SARS-CoV-2 transmission to humans. Care  
328 should be taken to communicate this, especially given the potential consequences of  
329 miscommunication for wildlife conservation. The bat research community in particular has  
330 expressed concern that negative framing of bats as the source of SARS-CoV-2 will impact public  
331 and governmental attitudes toward bat conservation<sup>71</sup>. In zoonotic virus research on bats, studies  
332 often over-emphasize human disease risks<sup>72</sup> and rarely mention ecosystem services provided by  
333 these animals<sup>73</sup>. Skewed communication can fuel negative responses against bats, including  
334 indiscriminate culling (i.e., reduction of populations by selective slaughter)<sup>74</sup>, which has already  
335 occurred in response to COVID-19 even outside of Asia (where spillover occurred)<sup>75</sup>.

336  
337 To minimize potential unintended negative impacts for bat conservation, public health and  
338 conservation responses should act in accordance with substantial evidence suggesting that  
339 culling has numerous negative consequences, not only threatening population viability of  
340 threatened bat species in shared roosts<sup>76</sup> but also possibly increasing viral transmission within  
341 the very species that are targeted<sup>77,78</sup>. Instead, bat conservation programs and long-term

342 ecological studies are necessary to help researchers understand viral ecology and find  
343 sustainable solutions for humans to live safely with wildlife. From another perspective, policy  
344 solutions aimed at limiting human-animal contact could potentially prevent virus establishment  
345 in novel species (e.g., as observed in mink farms<sup>51</sup>), especially in wildlife that may already face  
346 conservation challenges (e.g., North American bats threatened by an emerging disease, white-  
347 nose syndrome<sup>74,79</sup>). At least four bat species with confirmed white-nose syndrome symptoms or  
348 that can be infected by the fungal pathogen (*Eptesicus fuscus*, *Myotis lucifugus*, *M.*  
349 *septentrionalis*, *Tadarida brasiliensis*) are in our list of the 291 bat species most likely to be  
350 betacoronavirus hosts, and both *Myotis* species have already been heavily impacted by this fungal  
351 epidemic with over 90% reductions in their populations<sup>80</sup>.

352  
353 Substantial investments are already being planned to trace the wildlife origins of SARS-CoV-2.  
354 However, the intermediate progenitor virus may never be isolated from samples  
355 contemporaneous with spillover, and it may no longer be circulating in wildlife. MERS-CoV  
356 circulates continuously in camels<sup>81</sup> and SARS-CoV persisted in civets long enough to seed  
357 secondary outbreaks<sup>49,50</sup>, but the limited description of Pangolin-CoV symptoms suggests high  
358 mortality, potentially indicating a more transient epizootic such as Ebola die-offs in red river hogs  
359 (*Potamochoerus porcus*)<sup>18</sup>. In lieu of concrete data, our study provides no additional evidence  
360 implicating any particular species—or any particular pathway of spillover (e.g., wildlife trade,  
361 consumption of hunted animals)—as more or less likely. No specific scenario can be confirmed  
362 or rigorously interrogated by ecological models, and we explicitly warn against misinterpretation  
363 or misuse of our findings as evidence for adjacent policy decisions. Although policies that focus  
364 on particular potential reservoir species or target human-wildlife contact could reduce future  
365 spillovers, they will have a negligible bearing on the ongoing pandemic, as SARS-CoV-2 is highly  
366 transmissible within humans (e.g., unlike MERS-CoV or other zoonoses that are sustained in  
367 people by constant reintroduction). SARS-CoV-2 is likely to remain circulating in human  
368 populations until a vaccine is developed, regardless of immediate actions regarding wildlife.  
369 COVID-19 response must be informed by the best consensus evidence available and prioritize  
370 solutions that address immediate reduction of transmission through public health and policy  
371 channels. Meanwhile, we hope our proposed wildlife sampling priorities will help increase the  
372 odds of preventing the future emergence of novel betacoronaviruses.

373

374

## 375 **Acknowledgements**

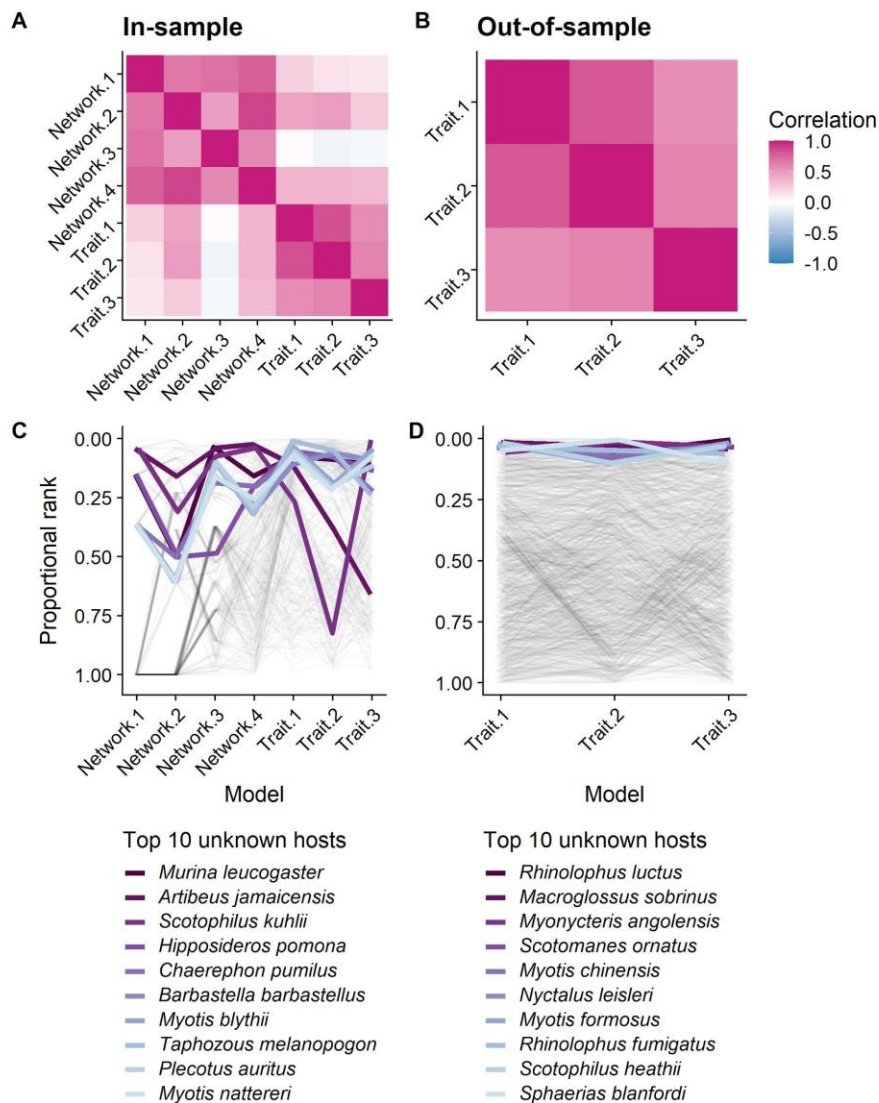
376 We thank Heather Wells for generously sharing thoughtful comments and code. The VERENA  
377 consortium is supported by L'Institut de Valorisation de Données (IVADO) through Université de  
378 Montreal. DJB was supported by an appointment to the Intelligence Community Postdoctoral  
379 Research Fellowship Program at Indiana University, administered by Oak Ridge Institute for  
380 Science and Education through an interagency agreement between the U.S. Department of Energy  
381 and the Office of the Director of National Intelligence.

382

383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395

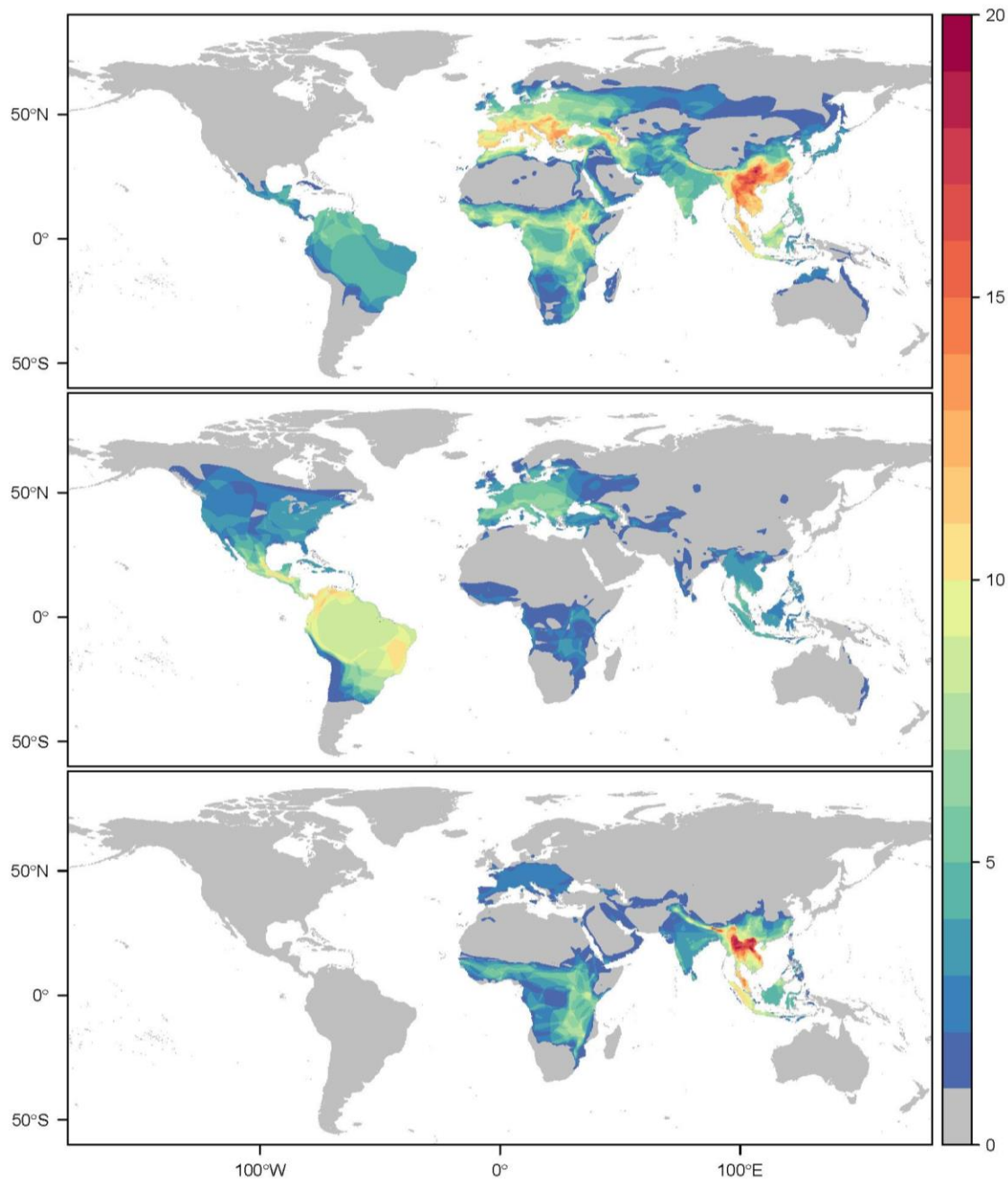
## Figures

**Figure 1. An ensemble of predictive models facilitates identification of likely betacoronavirus bat hosts.** The pairwise Spearman's rank correlations between models' ranked species-level predictions were generally substantial and positive (A,B). Models are arranged in decreasing order of their mean correlation with other models. In-sample predictions, expressed as host species' proportional rank (0 is the most likely host from a given model, 1 is the least likely host), varied significantly due to the uncertainty of network approaches (C). In contrast, species' proportional ranks were tightly correlated across out-of-sample predictive approaches, which relied on species traits (D). Each line represents a different bat species' proportional rank across models. The ten species with the highest mean proportional ranks across all models are highlighted in shades of purple.



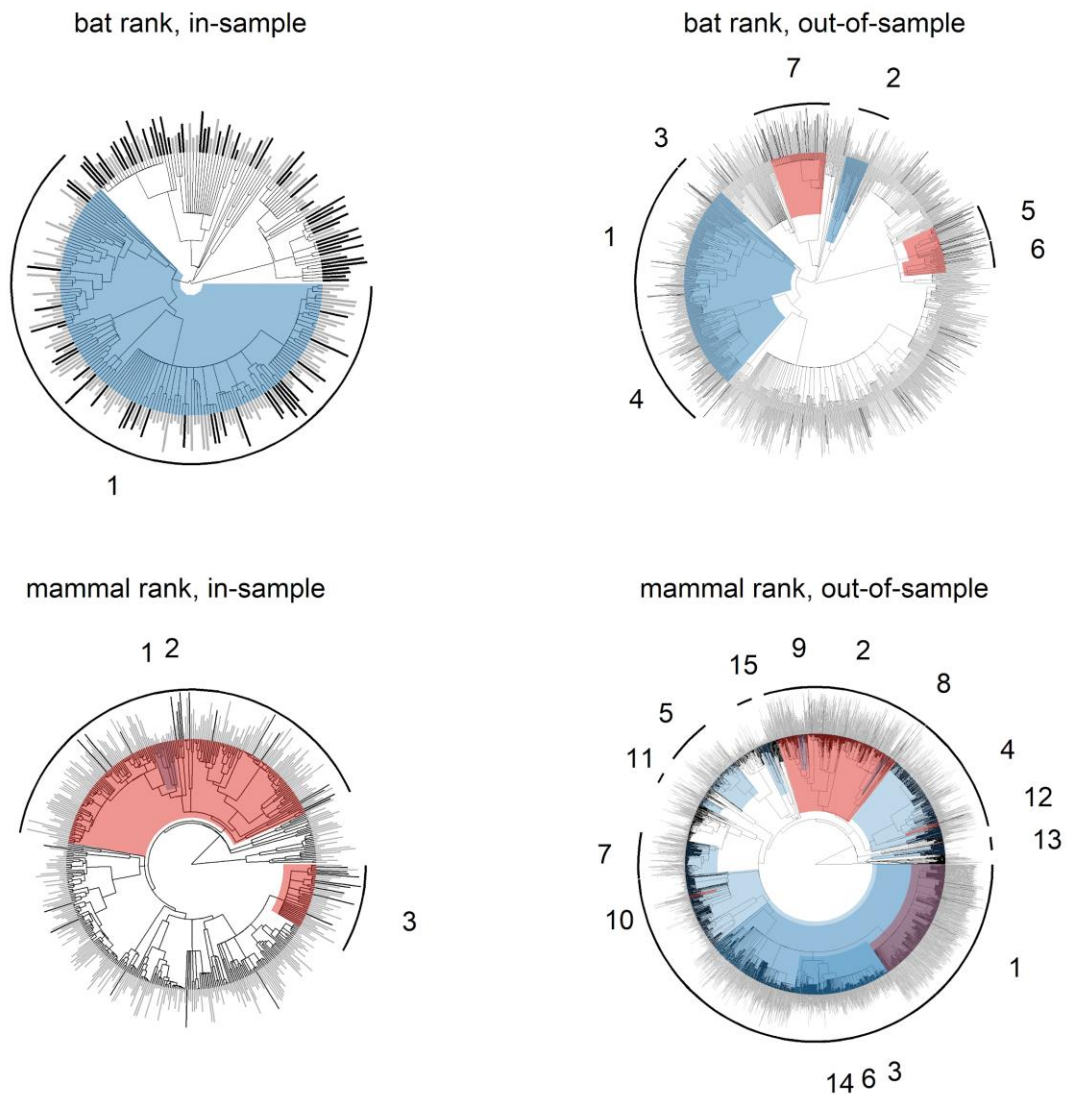
396

397 **Figure 2. Species richness of known and suspected betacoronavirus bat hosts.** Known hosts of  
398 betacoronaviruses (*top*) are found worldwide, but particularly in southern Asia and southern  
399 Europe. The top 50 predicted bat hosts with viral association records (*middle*) are mostly  
400 Neotropical, including several species of vampire bats. In contrast, the top 50 *de novo* bat host  
401 predictions based on phylogeny and ecological traits (*bottom*) are mostly clustered in Myanmar,  
402 Vietnam, and southern China, with none in the Neotropics or North America.  
403



404

405 **Figure 3. Phylogenetic distribution of predicted bat and mammal hosts of betacoronaviruses.**  
406 Bar height indicates mean predicted rank across the model ensemble (higher values = lower  
407 proportional rank score, more likely to be a host) and black indicates known betacoronavirus  
408 hosts. Colored regions indicate clades identified by phylogenetic factorization as significantly  
409 different in their predicted rank compared to the paraphyletic remainder; those clades more  
410 likely to contain a host are shown in red, whereas those less likely to contain a host are shown  
411 in blue. Results are displayed for bats and all mammals separately, stratified by in- and out-of-  
412 sample predictions. Numbers reference clade names, species richness, and mean predicted  
413 ranks as described in Extended Data Tables 1 and 2.  
414



415



416

417

## Methods.

418

419 The underlying conceptual aim of this study was to produce and synthesize several different  
420 models that predict and rank candidate reservoir species—each with different methods,  
421 assumptions, and framings—and to rapidly synthesize these into a consensus list. We broadly  
422 structured our study around two modeling targets: (1) produce rankings of likely bat hosts of  
423 betacoronaviruses and (2) identify potential non-bat mammal hosts. We developed a novel  
424 dataset that merged existing knowledge about the broader mammal-virus network with targeted  
425 data collection about coronaviruses; implemented seven modeling methods; synthesized these  
426 into an ensemble; and post-hoc identified taxonomic patterns in prediction using phylogenetic  
427 factorization.

428

### Host-Virus Association Data

429

430  
431 Entries were downloaded from GenBank on March 27th 2020 using the following search terms:  
432 Coronavirus, Coronaviridae, Orthocoronavirinae Alphacoronavirus, Betacoronavirus,  
433 Gammacoronavirus, and Deltacoronavirus. Data were sorted using a Python script that saved all  
434 available metadata regarding accession number, division, submission date, entry title, organism,  
435 genus, genome length, host classification, country, collection date, PubMed ID, journal containing  
436 associated publication, publication year, genome completeness, and the gene sequenced. The  
437 dataset was cleaned to remove duplicate entries, using GenBank accession number, and entries  
438 that did not correspond to viral sequences, using GenBank division. After cleaning, 31,473 entries  
439 remained, of which 25,628 had metadata regarding host species.

440

441 Data from GenBank were merged with the Host-Pathogen Phylogeny Project (HP3) dataset<sup>30</sup>. The  
442 HP3 dataset consists of 2,805 associations between 754 mammal hosts and 586 virus species,  
443 compiled from the International Committee on Taxonomy of Viruses (ICTV) database, and  
444 manually cleaned over a period of five years. Data collection on HP3 began in 2010 and has been  
445 static since 2017, but it still represents the most complete dataset on the mammal virome  
446 published with a high standard of data documentation. Several recent studies have used the HP3  
447 dataset to produce statistical models of viral sharing or zoonotic potential<sup>29,48,82</sup>, making it a  
448 comparable reference for a multi-model ensemble study.

449

450 Because of naming inconsistencies both within GenBank and between the two datasets (HP3 and  
451 GenBank), we used a two-step pipeline for taxonomic reconciliation. Viral names were matched  
452 to the ICTV 2019 master species list, up to the sub-genus level. Host species names were  
453 matched against GBIF using their species API with an automated Julia script, and processed to a  
454 fully cleaned set of names. This led to an harmonized dataset representing a global list of  
455 mammal-virus associations, from which the bat-coronavirus data can be extracted for  
456 downstream and specific modeling efforts. Because the HP3 dataset used an older version of the  
457 ICTV master list, and because not all host names in the GenBank metadata could be matched by



458 the GBIF species API (or could be solved unambiguously to the species level), some host-virus  
459 interactions were lost; this reinforces the need to careful data curation of taxonomic metadata if  
460 they are to enable and support predictive pipelines.

461

462

## 463 **Predictor Data**

464

### 465 *Phylogeny*

466

467 We used a supertree of extant mammals to unify modeling approaches incorporating host  
468 phylogeny<sup>31</sup>. Although more recent mammal supertrees exist, we used this particular phylogeny  
469 for consistency with trait datasets and several of the modeling frameworks included in our  
470 ensemble. We manually matched select bat species names between our edge list and this  
471 particular phylogeny. This included reverting any *Dermanura* to their former *Artibeus* designation  
472 (i.e., *D. phaeotis*, *D. cinerea*, *D. tolteca*)<sup>83</sup>, switching *Tadarida* species to either *Mops* or *Chaerephon*  
473 species (i.e., *Tadarida condylura* to *Mops condylurus*, *Tadarida plicata* to *Chaerephon plicatus*,  
474 *Tadarida pumila* to *Chaerephon pumilus*)<sup>84</sup>, and renaming *Myotis pilosus* to the more recent *Myotis*  
475 *ricketti*. *Chaerephon pusillus* was considered its own species but is now synonymous with  
476 *Chaerephon pumilus*<sup>84</sup>. Minor discrepancies between virus data and our phylogeny were also  
477 corrected (*Hipposideros commersonii* to *Hipposideros commersoni* [although more recently  
478 changed to *Macronycteris commersoni*], *Rhinolophus hildebrandti* to *Rhinolophus hildebrandtii*,  
479 *Neoromicia nana* to *Neoromicia nanus*). In other cases, some recently revised genera in our edge  
480 list were modified to match former genera in the mammal supertree: *Parastrellus hesperus* to  
481 *Pipistrellus hesperus*, and *Perimyotis subflavus* to *Pipistrellus subflavus*<sup>85</sup>. Lastly, some names in  
482 our edge list missing from the mammal supertree represent former subspecies being raised to  
483 full species rank, and names were reverted accordingly: *Artibeus planirostris* to *Artibeus*  
484 *jamaicensis*, *Miniopterus fuliginosus* to *Miniopterus schreibersii*, *Triaenops afer* to *Triaenops*  
485 *persicus*, and *Carollia sowelli* to *Carollia brevicauda*. Although we recognize that these are each  
486 now recognized as distinct species, in all cases our synonymized names are thought to be either  
487 sister taxa or very closely related.

488

### 489 *Ecological traits*

490

491 We used a previously published dataset of 63 ecological traits describing the morphology, life  
492 history, biogeography, and diet of 1,116 bat species. These data are drawn from a combination  
493 of PanTHERIA<sup>32</sup>, EltonTraits<sup>33</sup>, and the IUCN Red List range maps, and were previously cleaned in  
494 a study producing predictions of bat reservoirs of filoviruses<sup>27</sup>. Four redundant variables (two for  
495 human population density, mean potential evapotranspiration in range, and body mass) were  
496 eliminated prior to analyses, favoring variables with higher completeness.

497

### 498 *Correction for sampling bias*

499

500 To correct for sampling bias, in the style of several previous studies<sup>30,82</sup>, we used the number of  
501 peer-reviewed citations available on a given host as a measure of scientific sampling effort. We  
502 used the R package *easyPubMed* to scrape the number of citations in PubMed returned when  
503 searching each of the 1,116 bat names in the trait data on April 10, 2020.

504

## 505 **Modeling Approaches**

506

507 Our team produced an ensemble of seven statistical models (ED Tables 5 and 6), and applied  
508 them to generate a predictive set of seven models for bats and five for other mammals. Four use  
509 a network-theoretic component (k-nearest neighbors, linear filtering, trait-free plug-and-play, and  
510 scaled phylogeny), while three primarily used ecological traits as predictors (boosted regression  
511 trees, Bayesian additive regression trees, and neutral phylogeographic).

512

513 All eight approaches were used to generate predictions about potential bat hosts of  
514 betacoronaviruses. A subset of six were used to recommend potential non-bat mammal hosts of  
515 betacoronaviruses (k-nearest neighbor, linear filtering, scaled phylogeny, trait-free plug-and-play,  
516 and neutral phylogeographic). We did not use trait-based models to predict non-bat hosts,  
517 because assigning pseudoabsences to the vast majority (~3500 or more) of mammal species  
518 would likely lead to largely uninformative predictions, weighed against the 109 known  
519 betacoronavirus hosts (79 bats and 30 other mammals).

520

### 521 *Network model 1: k-Nearest Neighbors recommender*

522

523 We follow the methodology previously developed for the recommendation of species feeding  
524 interactions<sup>86</sup>. This method builds a recommender system internally based on the *k*-NN algorithm,  
525 under which candidate hosts are recommended for a virus from a pool constituted by the hosts  
526 of the *k* viruses with which it has the greatest overlap. Overlap (host sharing) is measured using  
527 Tanimoto similarity, which is the cardinality of the intersection of two sets divided by the  
528 cardinality of their union. To obtain the pairwise similarity between two viruses, this divides the  
529 number of shared hosts by the cumulative number of hosts. The *k* nearest neighbors of a virus  
530 are the *k* other viruses with which it has the highest Tanimoto similarity.

531

532 Hosts are then recommended by counting how many times they appear in these *k* neighbors, a  
533 quantity that ranges from 1 to *k*. We can impose arbitrary cutoffs by limiting the  
534 recommendations to the hosts that occur in at least *k*, *k*-1, etc, viruses. Previous leave-one-out  
535 validation of this model revealed that it is particularly effective for viruses with a reduced number  
536 of hosts, which is likely to be the case for emerging viruses. Furthermore, the performance of this  
537 model was not significantly improved by the addition of functional traits, making it acceptable to  
538 run on the association data only.

539

540 This model has been run two times; first, by measuring the similarity of viruses, and  
541 recommending hosts; second, by measuring the similarity of hosts, and recommending viruses.  
542 In all cases, only results for betacoronaviruses are reported.

543  
544 The outcome of this model should be subject to caution, as leave-one-out validation revealed that  
545 the success rate (*i.e.* ability to recover one interaction that has been removed) remained lower  
546 than 50% even when using  $k=8$ , and dropped as low as 5% when using  $k=1$  (the nearest-neighbor  
547 algorithm). This strongly suggests that the dataset of reported host-virus associations is  
548 extremely incomplete; therefore, the identification of the nearest neighbors can be biased by  
549 under-reported interactions, and this can result in noise in the prediction. This noise can be  
550 particularly important when the kNN technique operates on viruses, of which the bat dataset has  
551 only 15.

552  
553 *Network-based model 2: Linear filter recommender*

554  
555 Following Stock *et al.*<sup>87</sup>, we used a previously developed linear filter to infer potential missing  
556 interactions. This recommender system assumes that networks tend to be self-similar, and use  
557 this information to generate a score for an un-observed interaction that is a linear combination of  
558 the status of the interaction (relative weight of 1/4), relative degree of host and virus, and of the  
559 observed connectance of the network (all with relative weights of 1); as we are concerned with  
560 ranking interactions as opposed to examining the absolute value of the score, the penalization  
561 coefficient associated to the interaction being presumed absent could be omitted with no change  
562 in the ranking, but has been set to a low value instead. The scores returned by the linear filter are  
563 not directly related to the probability of the interaction existing in this context, but higher scores  
564 still indicate interactions that are more likely to exist. Indeed, known hosts of betacoronavirus  
565 typically scored higher.

566  
567 We used the zero-one-out approach to assess the performance of this model on the entire  
568 datasets. In all cases, non-interactions ranked lower than positive interactions even when entirely  
569 removing the penalization coefficient from the linear filter parameters, which suggests that the  
570 network structure (degree and connectance) is capturing a lot of information as to which species  
571 can interact. Note that as opposed to the k-NN method outlined above, the linear filter is  
572 symmetrical, *i.e.* it captures the properties of both host and virus at once.

573  
574 *Network-based model 3: Plug and play*

575  
576 For network problems, the “plug and play” model is a statistical approach that formulates Bayes’  
577 theorem for link prediction around the conditional density of traits of known associations  
578 compared to traits of every possible association in a network. The conditional density function is  
579 measured by using non-parametric kernel density estimators (implemented with the R package  
580 *np*), and the conditional ratio between them is used to estimate link “suitability”, a scale-free ratio.  
581 Compared to other machine learning methods that fit to training data iteratively, plug and play is

582 comparatively simple, and directly infers the most likely extensions of observed patterns in data.  
583 The plug and play was originally developed to forecast missing links in host-parasite networks<sup>36</sup>,  
584 but has since been used to model species distributions<sup>88</sup> and predict the global spread of human  
585 infectious diseases<sup>89</sup>. We used this model here to estimate suitability of host-virus interactions  
586 by first modeling the entire estimated network of host-virus interaction suitability, and ranking  
587 hosts that are not infected by betacoronaviruses by their estimated suitability for  
588 betacoronaviruses.

589  
590 The “plug and play” model is trained using either matched pairs of host and pathogen ecological,  
591 morphological, or phylogenetic traits<sup>36</sup>, or by using a latent approach<sup>89</sup> which considers the mean  
592 similarity of pathogens in their host ranges and the mean similarity of hosts in their pathogen  
593 communities as ‘traits’. We decided to use the latent approach, as host trait data was far more  
594 available than viral trait data. Further, the taxonomic scale considered for host (species) and virus  
595 (genus) differed, making the resolution of potential trait data different enough to potentially  
596 confound trait-based approaches in this modeling framework.

597  
598 Relative suitability of a host-virus association, as estimated by the “plug and play” model, is  
599 formulated as a density ratio estimation problem. The suitability of a host-virus association is  
600 quantified as the quotient of the distribution of latent trait values when an association was  
601 recorded over the distribution of all the latent trait values. As an attempt to control for sampling  
602 effort of mammal and bat host species, we included PubMed citation counts for host species (as  
603 described above) in the estimation of host-virus suitability. We explored host-pathogen suitability  
604 using the entire mammal-virus associations dataset, to maximize the available information on  
605 the network’s structure, and ranked host-pathogen pairs by their relative suitability value. From  
606 the final predictions, we subset out bat-specific predictions. When predicting, we set citation  
607 counts to the mean of training data, as a sampling bias correction.

608  
609 *Network-based model 4: Scaled-phylogeny*

610  
611 We apply the network-based conditional model of Elmasri *et al.*<sup>90</sup> for predicting missing links in  
612 bipartite ecological networks. The full model combines a hierarchical Bayesian latent score  
613 framework which accounts for the number of interactions per taxon, and a dependency among  
614 hosts based on evolutionary distances. To predict links based on evolutionary distance, the  
615 probability of a host-parasite interaction is taken as the sum of evolutionary distances to the  
616 documented hosts of that parasite. This allocates higher probabilities when a few closely related  
617 hosts, or many distantly related hosts interact with a parasite. In this way phylogenetic distances  
618 are combined with individual affinity parameters per taxa to model the conditional probability of  
619 an interaction.

620  
621 In ecological studies, it is common to use time-scaled phylogenies to quantify evolutionary  
622 distance among species<sup>91</sup>. We may use these fixed evolutionary distances for link prediction, but  
623 parasite taxa are known to be more or less constrained by phylogenetic distances among hosts<sup>92</sup>.

624 Further, phylogenies are hypotheses about evolutionary relationships and have uncertainties in  
625 the topology and relative distances among species<sup>93</sup>. Rather than treating phylogenetic distances  
626 as fixed, Elmasri *et al.*<sup>90</sup> re-scale the phylogeny by applying a macroevolutionary model of trait  
627 evolution. While any evolutionary model that re-scales the covariance matrix may be used, we use  
628 the early-burst model, which allows evolutionary change to accelerate or decelerate through  
629 time<sup>94</sup>. This different emphasis to be placed on deep versus recent host divergences when  
630 predicting links.

631  
632 We apply the model to a network of associations among host species and viral genera, and the  
633 mammal supertree, which allows us to leverage information from across the network to predict  
634 undocumented bat-betacoronavirus associations. We fit sets of models, applying both the full  
635 model, and the phylogeny-only model to both the bat-viral genera associations, and the mammal-  
636 viral genera associations. For each data-model combination we fit the model using ten-fold cross-  
637 validation holding out links for which there is a minimum of two observed interactions. The  
638 posterior interaction matrices resulting from each of the ten models are then averaged to  
639 generate predictions for all links in the network, with betacoronaviruses subset to generate the  
640 ensemble predictions.

641  
642 To assess predictive performance, we attempted to predict the held out interactions, and  
643 calculated AUC scores by thresholding predicted probabilities per fold, and taking an average  
644 across the 10 folds. In addition to AUC, we also assessed the model based on the percent of  
645 documented interactions accurately recovered. For the bat-viral genera data the full model  
646 resulted in an average AUC of 0.82 and recovered an average of 90.1% of held out interactions,  
647 while the phylogeny-only model showed increased AUC (0.86), but a decreased proportion of held-  
648 out interactions recovered (84.5%). Interestingly, the models for bat-virus genera associations  
649 had marginally worse predictive performance compared to the same models run on the larger  
650 network of mammal-virus associations (full model: AUC 0.88, 84.4% positive interactions  
651 recovered; phylogeny-only model:AUC: 0.88, 88.8% positive interactions recovered), indicating  
652 that predicting bat-betacoronavirus associations may benefit from including data on non-bat  
653 hosts. The models also estimated the scaling parameter ( $\eta$ ) of the early-burst model to be  
654 positive (average  $\eta=7.92$  for the full model run on the bat subset), indicating accelerating  
655 evolution compared to the input tree (ED Figure 6). This means that recent divergences are given  
656 more weight than deeper ones for determining bat-viral genera associations, which is consistent  
657 with recent work on viral sharing<sup>48,95</sup>.

658  
659 *Trait-based model 1: Boosted regression trees*

660  
661 Previous work has been highly successful in predicting zoonotic reservoirs using a combination  
662 of taxonomic, ecological, and geographic traits as predictors. This approach has been previously  
663 used to identify wildlife hosts of filoviruses<sup>27,96</sup>, flaviviruses<sup>28,97</sup>, henipaviruses<sup>23</sup>, *Borrelia*  
664 *burgdorferi*<sup>26</sup>, to predict mosquito vectors of flaviviruses<sup>98</sup>, and to predict rodent reservoirs and  
665 tick vectors of zoonotic viruses<sup>37,99</sup>. These approaches treat the presence of a specific virus (or

666 genus of viruses) or a zoonotic pathogen as an outcome variable, with negative values given for  
667 species not known to be hosts (pseudoabsences), and use machine learning to identify the  
668 characteristics that predispose animals to hosting pathogens of concern. By predicting the  
669 probability a given pseudoabsence is a false negative, the method can infer potential undetected  
670 or undiscovered host species.

671  
672 This approach has almost exclusively been implemented using boosted regression trees (BRT),  
673 a classification and regression tree (CART) machine learning method that became popular a  
674 decade ago for species distribution modeling.<sup>100</sup> Boosted regression trees develop an ensemble  
675 of classification trees which iteratively explain the residuals of previous trees, up to a fixed tree  
676 depth (usually between 3 and 5 splits). The incorporation of boosting allows the model, as it is fit,  
677 to progressively better explain poorly-fit cases within training data.

678  
679 We used boosted regression trees to identify trait profiles that predict bat hosts of  
680 betacoronaviruses, including all trait predictors from the trait database that met baseline  
681 coverage (< 50% missing values) and variation (< 97% homogenous) thresholds. For all model  
682 fitting, we specified a Bernoulli error distribution for our binary response variable and applied 10-  
683 fold cross validation to prevent overfitting (R package *gbm*). We started by fitting a global model  
684 to our full dataset, first specifying learning rate = 0.01 (shrinks the contribution of each tree to the  
685 model) and tree complexity = 4 (controls tree depth) as per default values and subsequently  
686 tuning to minimize cross validation error.

687  
688 We reduced the variable set by calling the *gbm.simp()* function, which computes and compares  
689 the mean change in cross validation error (deviance) produced by dropping different sets of least-  
690 contributing predictors. The final simplified model included 23 variables, plus citation counts,  
691 which we added to correct for sampling bias.

692  
693 We applied bootstrapping resampling methods to estimate uncertainty, using our tuned model to  
694 fit 1000 replicate models. For each model, training sets were assembled by randomly selecting  
695 with replacement 79 bat-coronavirus associations from the set of reported bat hosts and 79  
696 pseudoabsences. Trained models were used to generate relative influence coefficients for trait  
697 predictors and coronavirus host probabilities across all bat species. Partial dependence plots  
698 display relative influence coefficients and bootstrapped confidence intervals for the top ten  
699 contributing trait predictors. The medians of host probabilities were ranked and used to identify  
700 the top ten candidate host species. When predicting, we set citation counts to the mean of training  
701 data, as a sampling bias correction.

702  
703 *Trait based model 2: Bayesian additive regression trees*

704  
705 A similar workflow to trait-based model 1 was implemented using Bayesian additive regression  
706 trees (BART), an emerging machine learning tool that has similarities to more popular methods  
707 like random forests and boosted regression trees. BART adds several layers of methodological



708 innovation, and performs well in bakeoffs with other advanced machine learning methods.  
709 Several features make BART very convenient for modeling projects like these, including several  
710 easy-to-use implementations in R packages, built-in capacity to impute and predict on missing  
711 data, and easy construction of variable importance and partial dependence plots.

712  
713 Like other classification and regression tree methods, BART assigns the probability of a binary  
714 outcome variable by developing a set of classification trees - in this case, a sum-of-trees model -  
715 that split data ("branches") and assign values to terminal nodes ("leaves"). Whereas other similar  
716 methods generate uncertainty by adjusting data (e.g. random forests bootstrap training data and  
717 fit a tree to each bootstrap; boosted regression trees are usually implemented with iterated  
718 training-test splits to generate confidence intervals), BART generates uncertainty using an MCMC  
719 process. An initial sum-of-trees model is fit to the entire dataset, and then rulesets are adjusted  
720 in a limited and stochastic set of ways (e.g., adding a split; switching two internal nodes), with the  
721 sum-of-trees model backfit to each change. After a burn-in period, the cumulative set of sum-of-  
722 trees models is treated as a posterior distribution. This has some advantages over other methods,  
723 like boosted regression trees or random forests. In particular, posterior width directly measures  
724 model uncertainty (rather than approximating it by permuting training data), and a single model  
725 can be run (instead of an ensemble trained on smaller subsets of training data), allowing the  
726 model to use the full training dataset all at once.<sup>101</sup>

727  
728 Unlike many Bayesian machine learning methods, BART is easily implemented out-of-the-box, due  
729 to a limited set of customization needs. Three main priors control the fitting process: one usually-  
730 uniform prior on variable importance, one two-parameter negative power distribution on tree  
731 depth (preventing overfitting), and an inverse chi-squared distribution on residual variance. A set  
732 of well-performing priors from the original BART study<sup>102</sup> are widely used across R  
733 implementations for out-of-the-box settings, but can be further adjusted relative to modeling  
734 needs. In this study, we implemented BART models using a Dirichlet prior for variable importance  
735 (DART), a specification that is designed for situations with high dimensionality data that probably  
736 reflects a small number of true informative predictors. This often produces a much more reduced  
737 model without going through a stepwise variable selection process, which can be slow and very  
738 subject to stochasticity.<sup>101</sup>

739  
740 We implemented this approach using the *BART* package in R, using the bat-virus association  
741 dataset to generate an outcome variable, and the bat traits dataset as predictors. BART models  
742 were implemented with 200 trees and 10,000 posterior draws, using every trait feature that was  
743 at least 50% complete and < 97% homogenous (taken from TBM1).

744  
745 We tried four total implementations, based on two decisions: BART uncorrected and corrected  
746 for citation counts (BART-u, BART-c), and DART uncorrected and corrected for citation counts  
747 (DART-u, DART-c). All four models performed well, with little variation in predictive power  
748 measured by the area under the receiver operator curve calculated on training data (BART-u: AUC  
749 = 0.93; BART-c: AUC = 0.93; DART-u: AUC = 0.93; DART-c: 0.90; ED Figure 7). Across all models,

750 spatial variables had a high importance, including some regionalization (extent of range) and  
751 some variables capturing larger geographic range sizes, as did a diet of invertebrates (pulling out  
752 the phylogenetic signal of insectivorous bats; ED Figure 8).

753  
754 All models identified a number of “false negative” hosts that would be suitable based on a 10%  
755 false negative classification threshold for known betacoronavirus hosts (implemented with the R  
756 package ‘PresenceAbsence’). BART-u identified 217 missing hosts, BART-c identified 279  
757 missing hosts, DART-u identified 222 missing hosts, and DART-c identified 384 missing hosts,  
758 suggesting that this model most penalized overfitting as intended. As a result, we considered this  
759 model the most rigorous and powerful for inference, and used DART-c in the final model  
760 ensemble. We predicted across all 1,040 bats without recorded betacoronavirus associations,  
761 and ranked predicted probability. When predicting, we set citation counts to the mean of training  
762 data, as a sampling bias correction.

763  
764 *Trait based model 3: Phylogeographic neutral model*

765  
766 We used a previously published pairwise viral sharing model<sup>48</sup> to predict potential  
767 betacoronavirus hosts based on the sharing patterns of known hosts in a published dataset<sup>30</sup>.  
768 We used a generalised additive mixed model (GAMM), which was fitted in the first half of 2019  
769 using the *mgcv* package, with pairwise binary viral sharing (0/1 denoting if a species shares at  
770 least one virus) as a response variable. Explanatory variables include pairwise proportional  
771 phylogenetic distance and geographic range overlap (taken from the IUCN species ranges), with  
772 a multi-membership random effect to control for species-level sampling biases. The model was  
773 then used to predict the probability that a given species pair share at least one virus across 4196  
774 placental mammals with available data, producing a predicted viral sharing network that  
775 recapitulates a number of known macroecological patterns, as well as predicting reservoir hosts  
776 with surprising accuracy<sup>48</sup>. Subsetting this predicted sharing matrix, we listed the rank order of  
777 hosts most likely to share with all known betacoronavirus hosts in our datasets.

778  
779 *Rhinolophus-specific implementation of Trait-based model 3*

780  
781 We then repeated this process with sharing patterns of *Rhinolophus affinis* and *R. malayanus*  
782 specifically. Given the strong phylogenetic effect, the top 139 predictions were bat species:  
783 predominantly rhinolophids and hipposiderids. The top 20 predictions for both *R. malayanus* and  
784 *R. affinis* are displayed in ED Table 3 and 4. Notable predictions included the hog badger *Arctonyx*  
785 *collaris* (Carnivora: Mustelidae), which was examined for SARS-CoV antibodies in 2003 and is  
786 reported in wildlife markets<sup>7,103</sup>; a selection of civet cats (Carnivora: Viverridae) including *Viverra*  
787 species; the binturong (*Arctitis binturong*); and the masked palm civet (*Paguma larvata*), the latter  
788 of which were implicated in the chain of emergence for SARS-CoV<sup>49,50</sup>; and pangolins (Pholidota:  
789 Manidae) including *Manis javanica* and *Manis pentadactyla*, which have been hypothesised to be  
790 part of the emergence chain for SARS-CoV-2<sup>18,19</sup>.

791

792 Alongside these high-ranked species-level predictions, we visually examined how predictions  
793 varied across all mammal orders and families using the whole dataset (ED Figure 5). Pangolins  
794 (Pholidota), treeshrews (Scandentia), carnivores (Carnivora), hedgehogs (Erinaceomorpha), and  
795 even-toed ungulates (Artiodactyla) had high mean predicted probabilities. Investigating family-  
796 level sharing probabilities revealed that civets (Viverridae) and mustelids (Mustelidae) were  
797 responsible for the high Carnivora probabilities, and mouse deer (Tragulidae) and bovids  
798 (Bovidae) were mainly responsible for high probabilities in the Artiodactyla (ED Figure 6).

799

## 800 **Consensus Methods and Recommendations**

801

### 802 *Combining and ranking predictions*

803

804 For seven models predicting bat hosts of betacoronaviruses, and five models predicting mammal  
805 hosts of betacoronaviruses, we combined predictions—generated using the same standardized  
806 data—into one standardized dataset. All mammal models were trained on data including bats, but  
807 predictions were subset to exclude bats to focus on likely intermediate hosts.

808

809 Each study's unique output—a non-intercomparable mix of different definitions of suitability or  
810 probability of association—were transformed into proportional rank, where lower rank indicates  
811 higher evidence for association out of the total number of hosts examined. By rescaling all results  
812 to proportional ranks between zero and one, we also allowed comparison of in-sample and out-  
813 of-sample predictions across all models. Proportional ranks were averaged across models to  
814 generate one standardized list of predictions. This absorbed much of the variation in model  
815 performance (ED Figure 1) and produced a set of rankings that performed well.

816

817 We elected not to withhold any “test” data to measure model performance, given that each  
818 method deployed in the ensemble has been independently and rigorously tested and validated in  
819 previous publications. Instead, to maximize the amount of available training data for every model,  
820 we used full datasets in each model and measured performance on the full training data.

821

822 For bats, the final ensemble of models spanned a large range of performance on the training data,  
823 measured by the area under the receiver operator curve (AUC; Network 1: 0.624; Network 2: 0.987;  
824 Network 3: 0.514; Network 4: 0.726; Trait 1: 0.850; Trait 2: 0.902; Trait 3: 0.762), indicating that it  
825 was possible to suitably detect differences in model performance on the full data. The total  
826 ensemble of proportional ranks performed medium well (AUC = 0.791). We used known  
827 betacoronavirus associations to threshold each model and the ensemble predictions based on a  
828 10% omission threshold (90% sensitivity), and we again found a wide range in the number of  
829 predicted undiscovered bat hosts of betacoronaviruses (Network 1: 162 species; Network 2: 1;  
830 Network 3: 111; Network 4: 44; Trait 1: 425; Trait 2: 384; Trait 3: 720; total ensemble: 291 species).  
831 Given concerns about mammal model performance and biological accuracy (see Main Text), we  
832 elected not to apply this exercise to mammal hosts at large.

833

834 To visualize the spatial distribution of predicted bat hosts, we used the IUCN Red List database  
835 of species geographic distributions. We took the top 50 ranked in-sample predictions and top 50  
836 ranked out-of-sample predictions and combined these range maps to visualize species richness  
837 of top predicted hosts (Figure 3).

838

#### 839 *Phylogenetic factorization of ensemble models*

840

841 We used phylogenetic factorization to flexibly identify taxonomic patterns in the consensus  
842 proportional rankings of likely hosts of SARS-CoV-2. Phylogenetic factorization is a graph-  
843 partitioning algorithm that iteratively partitions a phylogeny in a series of generalized linear  
844 models to identify clades at any taxonomic level (e.g., rather than *a priori* comparing strictly  
845 among genera or family) that differ in a trait of interest<sup>45</sup>. Using the mammal supertree, we used  
846 the *phylofactor* package to partition proportional rank as a Gaussian-distributed variable. We  
847 determined the number of significant phylogenetic factors using a Holm's sequentially rejective  
848 5% cutoff for the family-wise error rate. We applied this algorithm across our four final ensemble  
849 prediction datasets: in-sample bat ranks, out-of-sample bat ranks, in-sample mammal ranks, and  
850 out-of-sample mammal ranks.

851

852 Using network and trait-based models within-sample, we identified only one bat clade with  
853 substantially different consensus proportional rankings, the Yangochiroptera ( $\bar{x}=0.55$  compared  
854 to 0.42 for the remaining bat phylogeny, the Yinpterochiroptera). Out of sample, using only trait-  
855 based models, we instead identified seven bat clades with different propensities to include unlikely  
856 or likely bat hosts of betacoronaviruses. Subclades of the New World superfamily Noctilionoidea  
857 broadly had higher proportional ranks ( $\bar{x}=0.72$ ), indicating lower predicted probability of being  
858 hosts, as did the Emballanuridae ( $\bar{x}=0.77$ ). In contrast, several subfamilies of the Old World fruit  
859 bats (Pteropodidae), including the Rousettinae, Cynopterinae, and Eidolinaei, all had lower mean  
860 ranks ( $\bar{x}=0.27$ ). Our models also collectively identified the Rhinolophidae as having lower ranks  
861 ( $\bar{x}=0.36$ ).

862

863 Using network models within-sample across non-volant mammals, we identified four clades with  
864 different proportional ranks. The largest clade was the Laurasiatheria (Artiodactyla,  
865 Perissodactyla, Carnivora, Pholidota, Soricomorpha, and Erinaceomorpha), which had lower  
866 proportional ranks (higher risk;  $\bar{x}=0.55$ ). Nested within this clade, the Cetacea had greater  
867 proportional ranks ( $\bar{x}=0.89$ ), indicating lower risk. A large subclade of the Murinae (Old World rats  
868 and mice) also had lower ranks ( $\bar{x}=0.52$ ). Out of sample, using only the biogeographic viral sharing  
869 model, we instead identified 15 clades with different proportional ranks. The first clade identified  
870 large swaths of the Muridae as having higher risk ( $\bar{x}=0.38$ ) as well as the Laurasiatheria ( $\bar{x}=0.50$ ).  
871 Old World primates had weakly lower risk ( $\bar{x}=0.65$ ), as did the Scuridae ( $\bar{x}=0.67$ ). The Cetacea and  
872 Pinnipedia both had greater proportional ranks ( $\bar{x}=0.89$  and  $\bar{x}=0.71$ ). Old World porcupines  
873 (Hystricidae) and the Erinaceidae (Paraechinus, Hemiechinus, Mesechinus, Erinaceus, Atelerix)  
874 both had greater risk ( $\bar{x}=0.48$  and  $\bar{x}=0.39$ ), while the Afrosoricida had higher ranks ( $\bar{x}=0.97$ ).

875

876 To assess potential discrepancy between taxonomic patterns in model ensemble predictions and  
877 those of simply host betacoronavirus status itself, we ran a secondary phylogenetic factorization  
878 treating host status as a Bernoulli-distributed variable, with the same procedure applied to  
879 determine the number of significant phylogenetic factors. To assess sensitivity of taxonomic  
880 patterns to sampling effort, we ran phylogenetic factorization with and without square-root  
881 transformed PubMed citations per species as a weighting variable (ED Figure 9).

882  
883 Without accounting for study effort, phylogenetic factorization of betacoronavirus host status  
884 identified one significant clade across the bat phylogeny, the Yangochiroptera, as having fewer  
885 positive species (4.71%) than the paraphyletic remainder (12.12%). When accounting for study  
886 effort, however, the single clade identified by phylogenetic factorization changed, with a subclade  
887 of the family Pteropodidae (the Rousettinae) having a greater proportion of positive species  
888 (28.6%). For non-volant mammals, phylogenetic factorization identified only one clade, the family  
889 Camelidae, as having more positive species (75%) than the tree remainder (0.68%).

890

#### 891 *Phylogenetic factorization of Rhinolophidae virus sharing*

892

893 Because phylogenetic patterns in predictions from our viral sharing model could vary across other  
894 taxonomic scales beyond order and family, we also used phylogenetic factorization to more  
895 flexibly identify host clades with different propensities to share viruses with *R. affinis* and *R.*  
896 *malayanus*. We partitioned rank as a Gaussian-distributed variable and again determined the  
897 number of significant phylogenetic factors using Holm's sequentially rejective 5% cutoff.

898

899 Within the Chiroptera, we identified 10 clades with different propensities to share viruses with *R.*  
900 *affinis* and 5 clades with different propensities to share viruses with *R. malayanus*. For both bats,  
901 the top clade was the family Rhinolophidae, reinforcing phylogenetic components of the  
902 biogeographic model and highlighting the greater likelihood of viral sharing (mean rank  $\bar{x}$ =40 for  
903 *R. affinis*,  $\bar{x}$ =42 for *R. malayanus*). For *R. affinis*, several individual bat species had lower risks of  
904 viral sharing (e.g., *Myotis leibii*,  $\bar{x}$ =4100; *Pteropus insularis*,  $\bar{x}$ =3157; *Nyctimene aello*,  $\bar{x}$ =2497;  
905 *Chaerephon chapini*,  $\bar{x}$ =2497). The Megadermatidae, Nycteridae, and Hipposideridae (under which  
906 the PanTHERIA dataset includes the genus *Rhinonycteris*, although this is now considered a  
907 separate family, the Rhinonycteridae<sup>104</sup>) collectively had greater likelihood of viral sharing  
908 ( $\bar{x}$ =557), as did the Vespertilionidae ( $\bar{x}$ =704).

909

910 Across the non-volant mammals, we identified 7 clades with different propensities to share  
911 viruses with *R. affinis* and only 1 clade with different propensities to share viruses with *R.*  
912 *malayanus*. For both bat species, the first and primary clade was the Ferungulata (Artiodactyla,  
913 Perissodactyla, Carnivora, Pholidota, Soricomorpha, and Erinaceomorpha), which had lower ranks  
914 (higher viral sharing;  $\bar{x}$ =2084). For viral sharing with *R. affinis*, the Sciuridae was more likely to  
915 share viruses ( $\bar{x}$ =1948), as was the Scandentia ( $\bar{x}$ =1416) and many members of the Colobinae  
916 ( $\bar{x}$ =1958). However, members of the tribe Muntiacini (genera *Elaphodus* and *Muntiacus*) had  
917 especially high likelihoods of viral sharing and low rank ( $\bar{x}$ =361).

918

919 **Data and Code Availability**

920

921 The standardized data on betacoronavirus associations, and all associated predictor data, is  
922 available from the VERENA consortium's Github ([github.com/viralemergence/virionette](https://github.com/viralemergence/virionette)). All  
923 modeling teams contributed an individual repository with their methods, which are available in  
924 the organizational directory ([github.com/viralemergence](https://github.com/viralemergence)). All code for analysis, and a working  
925 reproduction of each authors' contributions, is available from the study repository  
926 ([github.com/viralemergence/Fresnel](https://github.com/viralemergence/Fresnel)).  
927



928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938

## Extended Data

**Extended Data Table 1. Results of phylogenetic factorization applied to predicted rank probabilities for bats.** The number of retained phylogenetic factors (following a 5% family-wise error rate applied to GLMs), taxa corresponding to those clades, number of species per clade, and mean predicted rank probabilities for the clade compared to the paraphyletic remainder are shown stratified by models applied in- and out-of-sample.

Sample	Factor	Taxa	Tips	Clade	Other
in	1	Yangochiroptera	160	0.549	0.422
out	1	Mystacinidae, Noctilionidae, Mormoopidae, Phyllostomidae	161	0.724	0.488
out	2	Mosia, Emballonura, Coleura, Rhynchonycteris, Cyttarops, Diclidurus, Centronycteris, Cormura, Saccopteryx, Balantiopteryx, Peropteryx	31	0.774	0.516
out	3	Thyropteridae, Furipteridae, Natalidae	12	0.853	0.520
out	4	Molossidae	98	0.595	0.517
out	5	Rousettus, Megaloglossus, Eidolon, Myonycteris, Plerotes, Casinycteris, Scotonycteris, Nanonycteris, Hypsignathus, Epomops, Micropteropus, Epomophorus	35	0.267	0.533
out	6	Sphaerias, Alionycteris, Otopteropus, Haplonycteris, Latidens, Penthetor, Thoopterus, Aethalops, Balionycteris, Chironax, Dyacopterus, Ptenochirus, Megaerops, Cynopterus	26	0.263	0.531
out	7	Rhinolophidae	73	0.360	0.536

939  
940

941 **Extended Data Table 2. Results of phylogenetic factorization applied to predicted rank**  
 942 **probabilities for all mammals.** The number of retained phylogenetic factors (following a 5%  
 943 family-wise error rate applied to GLMs), taxa corresponding to those clades, number of species  
 944 per clade, and mean predicted rank probabilities for the clade compared to the paraphyletic  
 945 remainder are shown stratified by models applied in- and out-of-sample.  
 946  
 947

Sample	Factor	Taxa	Tips	Clade	Other
in	1	Phocoenidae, Delphinidae, Tursiops, Monodontidae, Physeteridae, Balaenopteridae, Eschrichtiidae	12	0.889	0.611
in	2	Artiodactyla, Perissodactyla, Carnivora, Pholidota, Erinaceomorpha, Soricomorpha	173	0.549	0.661
in	3	Lophuromys, Micaelamys, Apodemus, Arvicanthis, Bandicota, Madromys, Dasymys, Hydromys, Lemniscornys, Mastomys, Mus, Pelomys, Niviventer, Otomys, Praomys, Rattus, Vandeleuria	38	0.520	0.627
out	1	Abditomys, Bullimus, Limnomys, Tarsomys, Tryphomys, Acomys, Lophuromys, Uranomys, Aethomys, Micaelamys, Anisomys, Chiruromys, Coccymys, Crossomys, Hyomys, Leptomys, Lorentzimys, Pseudohydromys, Paraleptomys, Macruromys, Mallomys, Microhydromys, Parahydromys, Pogonomelomys, Abeomelomys, Solomys, Xenuromys, Apodemus, Tokudaia, Apomys, Crunomys, Chrotomys, Rhynchomys, Arvicanthis, Bandicota, Batomys, Carpomys, Crateromys, Berylmys, Bunomys, Chiromyscus, Chiropodomys, Hapalomys, Haeromys, Colomys, Nilopegamys, Conilurus, Leporillus, Mesembriomys, Melomys, Protochromys, Mammelomys, Paramelomys, Uromys, Zyzomys, Leggadina, Notomys, Pseudomys, Mastacomys, Madromys, Cremnomys, Millardia, Dacnomys, Dasymys, Dephomys, Hybomys, Hydromys, Xeromys, Desmomys, Diomys, Diplothrix, Echiothrix, Margaretamys, Melasmothrix, Tateomys, Eropeplus, Lenomys, Golunda, Grammomys, Thallomys, Hadromys, Heimyscus, Hylomyscus, Komodomys, Papagomys, Oenomys, Thamnomys, Lemniscornys, Lenothrix, Leopoldamys, Malacomys, Praomys, Myomyscus, Mastomys, Maxomys, Micromys, Muriculus, Mus, Mylomys, Pelomys, Stenocephalemys, Nesokia, Niviventer, Otomys, Parotomys, Palawanomys, Paruromys, Phloeomys, Pithecheir, Pogonomys, Rattus, Rhabdomys, Srilankamys, Nesoromys, Stochomys, Sundamys, Taeromys, Vandeleuria, Vernaya, Zelotomys	510	0.382	0.672
out	2	Artiodactyla, Perissodactyla, Carnivora, Pholidota	505	0.495	0.651
out	3	Anomaluridae, Pedetidae, Dipodidae, Cricetidae, Muridae, Nesomyidae, Calomyscidae, Spalacidae, Platacanthomyidae	779	0.643	0.622
out	4	Talpidae, Erinaceomorpha, Soricidae	357	0.630	0.627
out	5	Cercopithecidae, Hominidae, Hylobatidae	139	0.649	0.626

948

949 **Extended Data Table 2, continued.** (Page 2 of 2)

950

Sample	Factor	Taxa	Tips	Clade	Other
out	6	Abrawayamys, Handleyomys, Aepeomys, Thomasomys, Abrothrix, Akodon, Necromys, Deltamys, Thaptomys, Andalgalomys, Auliscomys, Loxodontomys, Phyllotis, Paralomys, Graomys, Andinomys, Bibimys, Kunsia, Scapteromys, Blarinomys, Calomys, Chelemys, Chilomys, Chinchillula, Delomys, Eligmodontia, Euneomys, Galenomys, Geoxus, Holochilus, Landomys, Pseudoryzomys, Irenomys, Lenoxus, Melanomys, Microryzomys, Neacomys, Nectomys, Neotomys, Nesoryzomys, Notiomys, Oecomys, Oligoryzomys, Oryzomys, Oxymycterus, Brucepattersonius, Phaenomys, Podoxymys, Punomys, Reithrodon, Rhagomys, Rhipidomys, Scolomys, Sigmodontomys, Thalpomys, Wiedomys, Wilfredomys, Juliomys, Zygodontomys, Anotomys, Chibchanomys, Ichthyomys, Neusticomys, Rheomys, Sigmodon, Nyctomys, Otonyctomys, Ototylomys, Tylomys, Baiomys, Scotinomys, Ochrotomys, Habromys, Neotomodon, Podomys, Osgoodomys, Megadontomys, Peromyscus, Onychomys, Isthmomys, Reithrodontomys, Hodomys, Xenomys, Neotoma, Nelsonia	397	0.703	0.616
out	7	Tamiasciurus, Sciurus, Rheithrosciurus, Microsciurus, Syntheosciurus, Pteromys, Petaurista, Belomys, Biswamoyopterus, Trogopterus, Pteromyscus, Aeromys, Eupetaurus, Aeretes, Glaucomys, Eoglaucomys, Hylopetes, Petinomys, Petaurillus, Iomys, Ratufa, Callosciurus, Glyphotes, Lariscus, Menetes, Rhinosciurus, Funambulus, Tamiops, Dremomys, Exilisciurus, Hyosciurus, Prosciurillus, Rubrisciurus, Nannosciurus, Sundasciurus	139	0.672	0.625
out	8	Phocoenidae, Delphinidae, Tursiops, Monodontidae, Physeteridae, Balaenopteridae, Eschrichtiidae	12	0.889	0.626
out	9	Odobenidae, Otariidae, Phocidae	33	0.714	0.626
out	10	Hystriidae	11	0.482	0.627
out	11	Caprolagus, Poelagus, Lepus, Oryctolagus	33	0.642	0.627
out	12	Paraechinus, Hemiechinus, Mesechinus, Erinaceus, Atelerix	15	0.388	0.628
out	13	Afrosoricida	41	0.970	0.623
out	14	Castoridae, Heteromyidae, Geomyidae, Octodontidae, Ctenodactylidae, Ctenomyidae, Abrocomidae, Caviidae, Dinomyidae, Petromuridae, Dasypsectidae, Myocastoridae, Echimyidae, Erethizontidae, Capromyidae, Cuniculidae, Thryonomyidae, Bathyergidae, Chinchillidae	295	0.872	0.603
out	15	Cheirogaleidae, Indriidae, Daubentoniidae, Lemuridae, Lepilemuridae	48	0.921	0.623

951

952 **Extended Data Table 3. Predicted high-similarity bat hosts sharing with *Rhinolophus affinis* and**  
 953 ***R. malayanus*.** Species on these lists may be particularly likely to be the ultimate evolutionary  
 954 origin of SARS-CoV-2, or a closely-related virus prior to recombination in an intermediate host.  
 955 Predictions are made based just on the average viral sharing probability inferred for the two hosts  
 956 from the phylogeography model (Trait-based 3). (\* Note that the two species have high sharing  
 957 probabilities with each other, potentially indicating that efforts to trace the origins of SARS-CoV-  
 958 2 are already very close to their target.)

959  
 960

<b>Rhinolophus affinis</b>	<b>Rhinolophus malayanus</b>
1. <i>Rhinolophus macrotis</i> (P=0.84)	1. <i>Rhinolophus shameli</i> (P=0.87)
2. <i>Rhinolophus stheno</i> (P=0.83)	2. <i>Rhinolophus coelophyllus</i> (P=0.84)
3. <i>Rhinolophus malayanus</i> (P=0.82)	3. <i>Rhinolophus thomasi</i> (P=0.84)
4. <i>Rhinolophus acuminatus</i> (P=0.81)	4. <i>Rhinolophus affinis</i> (P=0.82)
5. <i>Rhinolophus pearsonii</i> (P=0.78)	5. <i>Rhinolophus marshalli</i> (P=0.82)
6. <i>Rhinolophus shameli</i> (P=0.78)	6. <i>Rhinolophus pearsonii</i> (P=0.82)
7. <i>Rhinolophus thomasi</i> (P=0.78)	7. <i>Rhinolophus yunanensis</i> (P=0.79)
8. <i>Rhinolophus sinicus</i> (P=0.77)	8. <i>Rhinolophus paradoxolophus</i> (P=0.78)
9. <i>Rhinolophus trifoliatus</i> (P=0.76)	9. <i>Rhinolophus macrotis</i> (P=0.76)
10. <i>Rhinolophus marshalli</i> (P=0.72)	10. <i>Rhinolophus acuminatus</i> (P=0.75)
11. <i>Rhinolophus shortridgei</i> (P=0.71)	11. <i>Rhinolophus siamensis</i> (P=0.75)
12. <i>Rhinolophus luctus</i> (P=0.7)	12. <i>Rhinolophus rouxii</i> (P=0.72)
13. <i>Rhinolophus sedulus</i> (P=0.7)	13. <i>Rhinolophus stheno</i> (P=0.71)
14. <i>Rhinolophus rouxii</i> (P=0.69)	14. <i>Rhinolophus luctus</i> (P=0.69)
15. <i>Rhinolophus pusillus</i> (P=0.68)	15. <i>Rhinolophus trifoliatus</i> (P=0.65)
16. <i>Rhinolophus ferrumequinum</i> (P=0.67)	16. <i>Rhinolophus pusillus</i> (P=0.62)
17. <i>Rhinolophus lepidus</i> (P=0.67)	17. <i>Rhinolophus borneensis</i> (P=0.6)
18. <i>Hipposideros pomona</i> (P=0.66)	18. <i>Hipposideros lylei</i> (P=0.59)
19. <i>Rhinolophus celebensis</i> (P=0.66)	19. <i>Rhinolophus shortridgei</i> (P=0.59)
20. <i>Rhinolophus paradoxolophus</i> (P=0.66)	20. <i>Rhinolophus sinicus</i> (P=0.59)

961

962 **Extended Data Table 4. Predicted high-similarity non-bat hosts sharing with *Rhinolophus affinis***  
 963 **and *R. malayanus*.** Species on these lists may be particularly suitable as stepping stones for  
 964 betacoronavirus transmission from bats into humans, including potentially for SARS-CoV-2 and  
 965 other SARS-like viruses. Predictions are made based just on the average viral sharing probability  
 966 inferred for the two hosts from the phylogeography model (Trait-based 3). Species' binomial  
 967 names are included alongside their families.  
 968

<i>Rhinolophus affinis</i>		<i>Rhinolophus malayanus</i>	
1. <i>Arctonyx collaris</i> (P=0.33)	Mustelidae	1. <i>Arctonyx collaris</i> (P=0.29)	Mustelidae
2. <i>Budorcas taxicolor</i> (P=0.33)	Bovidae	2. <i>Herpestes urva</i> (P=0.28)	Herpestidae
3. <i>Viverra zangalla</i> (P=0.32)	Viverridae	3. <i>Lutrogale perspicillata</i> (P=0.28)	Mustelidae
4. <i>Manis javanica</i> (P=0.3)	Manidae	4. <i>Melogale personata</i> (P=0.27)	Mustelidae
5. <i>Mustela altaica</i> (P=0.3)	Mustelidae	5. <i>Viverra megaspila</i> (P=0.26)	Viverridae
6. <i>Ursus thibetanus</i> (P=0.3)	Ursidae	6. <i>Arctictis binturong</i> (P=0.25)	Viverridae
7. <i>Cynogale bennettii</i> (P=0.29)	Viverridae	7. <i>Euroscaptor klossi</i> (P=0.25)	Talpidae
8. <i>Elaphodus cephalophus</i> (P=0.29)	Cervidae	8. <i>Lutra sumatrana</i> (P=0.25)	Mustelidae
9. <i>Lutrogale perspicillata</i> (P=0.29)	Mustelidae	9. <i>Sus scrofa</i> (P=0.25)	Suidae
10. <i>Viverricula indica</i> (P=0.29)	Viverridae	10. <i>Capricornis milneedwardsii</i> (P=0.23)	Bovidae
11. <i>Capricornis sumatraensis</i> (P=0.28)	Bovidae	11. <i>Manis javanica</i> (P=0.23)	Manidae
12. <i>Chimarrogale himalayica</i> (P=0.28)	Soricidae	12. <i>Manis pentadactyla</i> (P=0.23)	Manidae
13. <i>Helarctos malayanus</i> (P=0.28)	Ursidae	13. <i>Mustela nudipes</i> (P=0.23)	Mustelidae
14. <i>Herpestes javanicus</i> (P=0.27)	Herpestidae	14. <i>Paguma larvata</i> (P=0.23)	Viverridae
15. <i>Hylomys suillus</i> (P=0.27)	Erinaceidae	15. <i>Panthera pardus</i> (P=0.23)	Felidae
16. <i>Mustela kathiah</i> (P=0.27)	Mustelidae	16. <i>Viverra zibetha</i> (P=0.23)	Viverridae
17. <i>Capricornis milneedwardsii</i> (P=0.26)	Bovidae	17. <i>Bandicota savilei</i> (P=0.22)	Muridae
18. <i>Catopuma temminckii</i> (P=0.26)	Felidae	18. <i>Chrotogale owstoni</i> (P=0.22)	Viverridae
19. <i>Crocidura negligens</i> (P=0.26)	Soricidae	19. <i>Crocidura fuliginosa</i> (P=0.22)	Soricidae
20. <i>Capricornis thar</i> (P=0.25)	Bovidae	20. <i>Crocidura vorax</i> (P=0.22)	Soricidae

969

970 **Extended Data Table 5. Taxonomic scale of model training data and predictive**  
 971 **implementation.** Notes: (1) These models generated predictions of sharing with *Rhinolophus*  
 972 *affinis* over all non-human mammals in the HP3 dataset, then subsetted to bats. (2) In these  
 973 models, bat-betacoronavirus predictions are based on a subset of binary outcomes for known  
 974 association with betacoronaviruses, without any other viruses included.  
 975

Model approach	Training data scale	Bat <i>Betacoronavirus</i> predictions	Mammal-wide <i>Betacoronavirus</i> predictions
<b>Network-based 1</b> k-Nearest neighbors	Bat-virus	✓	
<b>Network-based 1</b> k-Nearest neighbors	Mammal-virus		✓
<b>Network-based 2</b> Linear filter	Bat-virus	✓	
<b>Network-based 2</b> Linear filter	Mammal-virus		✓
<b>Network-based 3</b> Plug-and-play	Mammal-virus <sup>1</sup>	✓	✓
<b>Network-based 4</b> Scaled-phylogeny	Bat-virus	✓	
<b>Network-based 4</b> Scaled-phylogeny	Mammal-virus		✓
<b>Trait-based 1</b> Boosted regression trees	Bat-betacoronavirus <sup>2</sup>	✓	
<b>Trait-based 2</b> Bayesian additive regression trees	Bat-betacoronavirus <sup>2</sup>	✓	
<b>Trait-based 3</b> Neutral phylogeographic	Mammal-virus <sup>1</sup>	✓	✓

976

977



978 **Extended Data Table 6. Data scale of prediction, by method.** Some methods use  
 979 pseudoabsences to expand the scale of prediction but still only analyze existing data, with no out-  
 980 of-sample inference, while others can predict freshly onto new data. (\* Training data from the HP3  
 981 database uses pseudoabsences, but no new ones are generated in this study that modify the  
 982 model or the bat-virus association dataset)

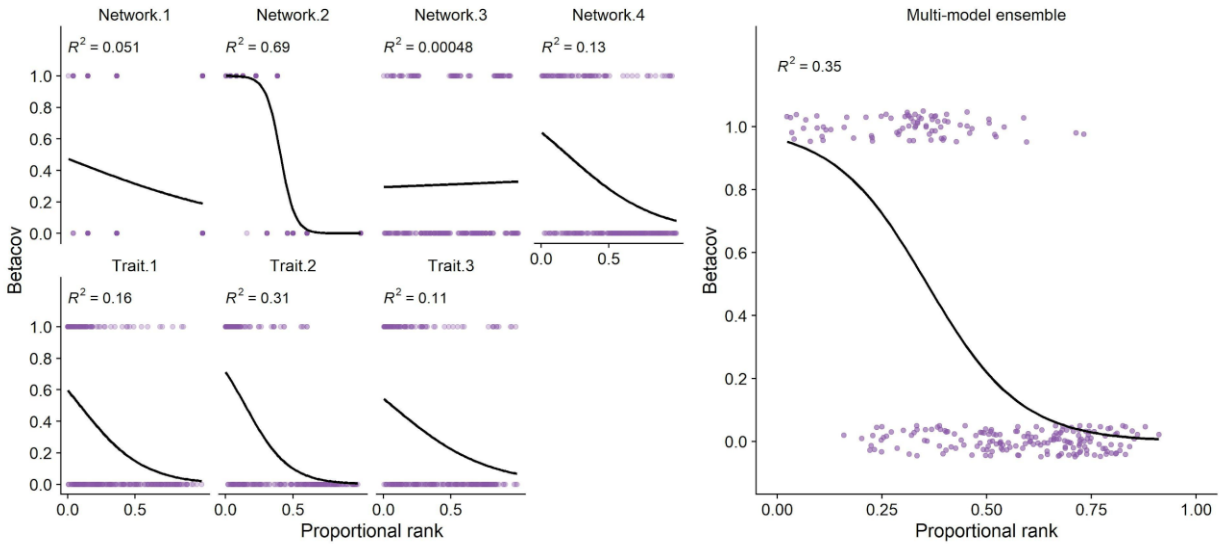
983

<b>Model approach</b>	<b>Prediction on hosts without known associations (out-of-sample)</b>	<b>Predictive extent and use of pseudoabsences</b>
<b>Network-based 1</b> k-Nearest neighbors	No	Only predicts link probabilities among species in the association data
<b>Network-based 2</b> Linear filter	No	Only predicts link probabilities among species in the association data
<b>Network-based 3</b> Plug-and-play	No	Uses pseudoabsences to predict over all mammals in association data, using latent approach
<b>Network-based 4</b> Scaled-phylogeny	No	Only predicts link probabilities among species in the association data
<b>Trait-based 1</b> Boosted regression trees	Yes	Uses pseudoabsences for all bats in trait data to predict over all species, including those without known associations
<b>Trait-based 2</b> Bayesian additive regression trees	Yes	Uses pseudoabsences for all bats in trait data to predict over all species, including those without known associations
<b>Trait-based 3</b> Neutral phylogeographic	Yes	Trains on a broader network, and predicts sharing probabilities among any mammals in phylogeny and IUCN range map data

984

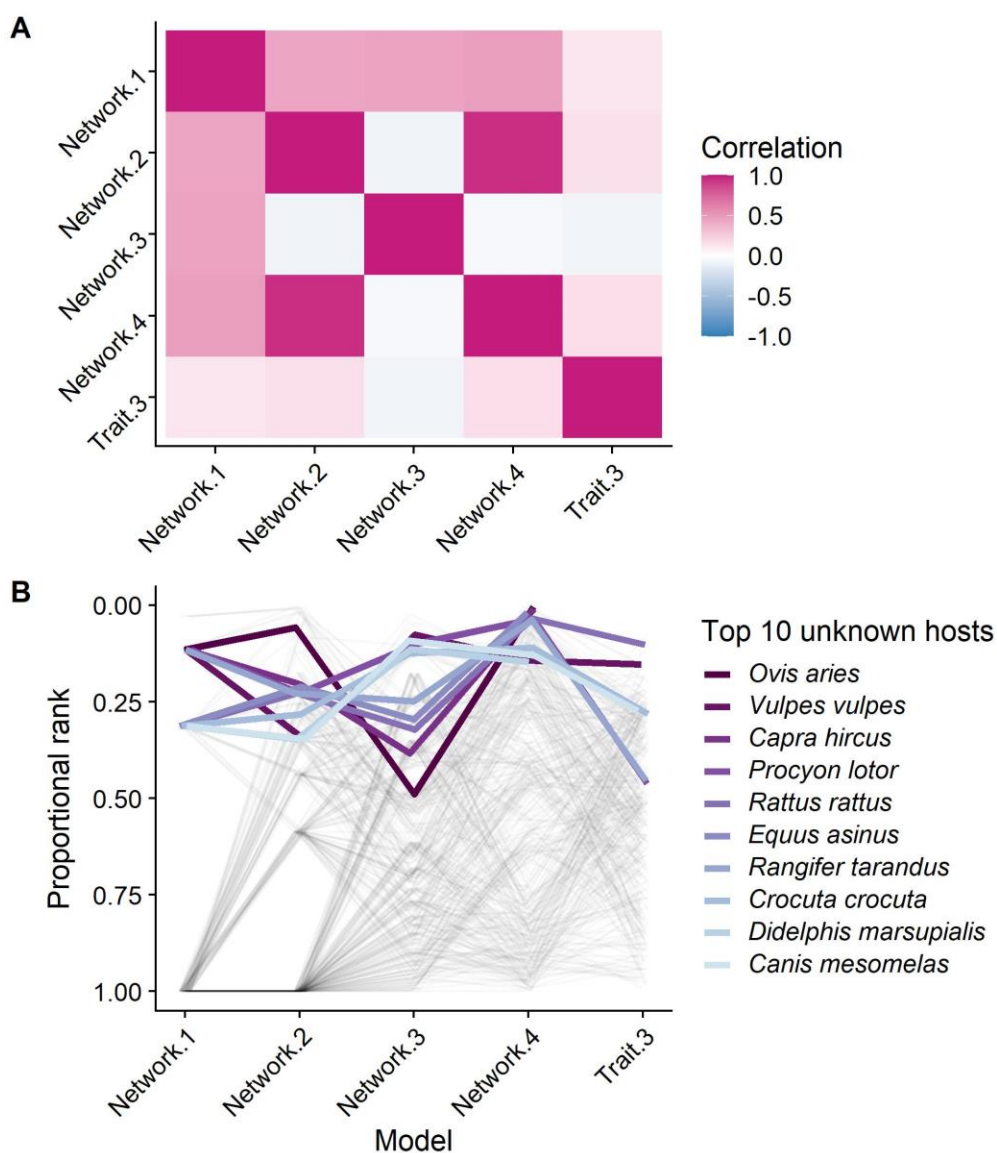
985

986 **Extended Data Figure 1. Bat models perform more strongly together than in isolation.** Curves  
987 show observed betacoronavirus hosts against predicted proportional ranks from seven individual  
988 models, and incorporated into one multi-model ensemble. Black lines show a binomial GLM fit to  
989 the predicted ranks against the recorded presence or absence of known betacoronavirus  
990 associations. Points are jittered to reduce overlap.  
991

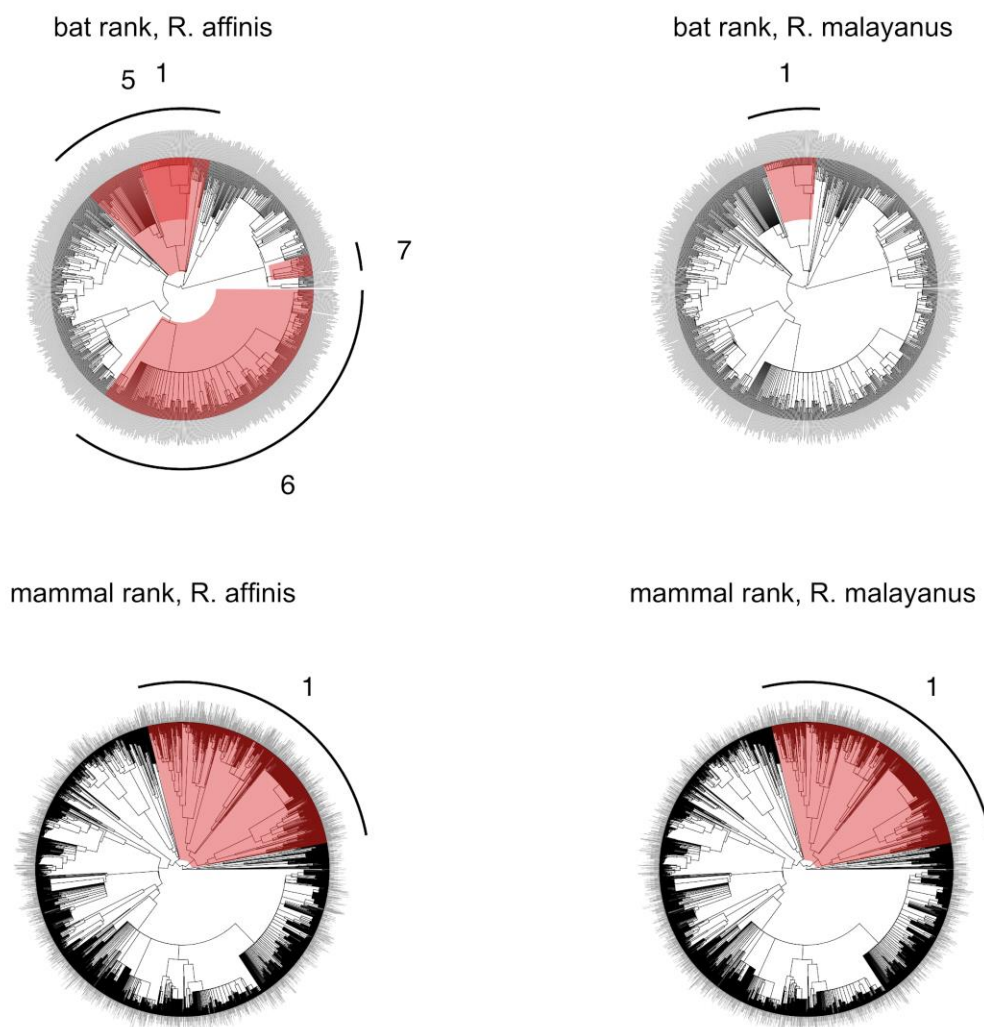


992

993 **Extended Data Figure 2. Poor concordance among predictive models for mammal hosts of**  
994 **betacoronaviruses.** The pairwise Spearman's rank correlations between models' ranked species-  
995 level predictions were generally low (A). In-sample predictions varied significantly and heavily  
996 prioritized domestic animals and well-studied hosts (B). The ten species with the highest mean  
997 proportional ranks across all models are highlighted in shades of purple. Only in-sample  
998 predictions are displayed because only one model (Trait-based 3) was able to predict out of  
999 sample for all mammals.

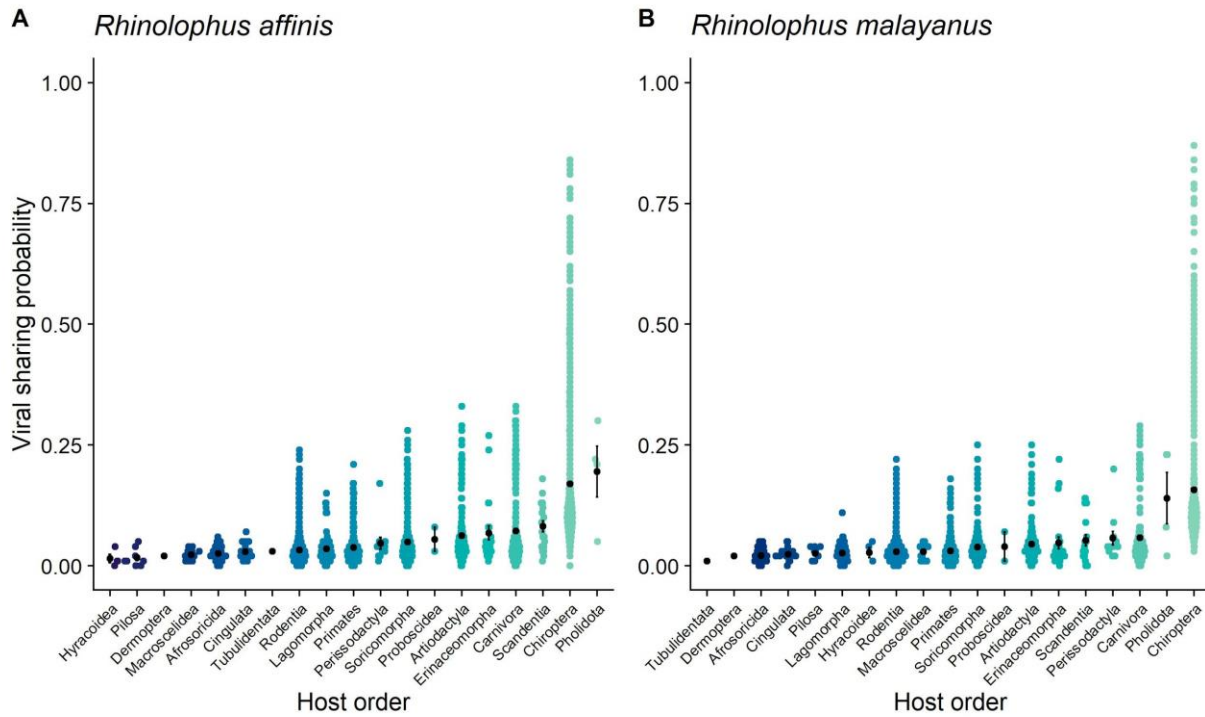


1001 **Extended Data Figure 3.** Results of phylogenetic factorization applied to predicted ranks of virus  
1002 sharing with *Rhinolophus affinis* and *Rhinolophus malayanus*. Colored regions indicate clades  
1003 identified as significantly different in their predicted rank compared to the paraphyletic remainder;  
1004 those more likely to share a virus with the *Rhinolophus* are shown in red, whereas those less likely  
1005 to share a virus are shown in blue. Bar height indicates predicted rank (higher values = lower rank  
1006 score, more likely share virus). Results are displayed for bats and remaining mammals separately.  
1007 Mammal-wide clades with high propensities to share viruses with *R. affinis* based solely on their  
1008 phylogeography included the treeshrews (Scandentia), Old World monkeys (Colobinae), and both  
1009 tufted and barking deer (Muntiacini).  
1010  
1011



1012

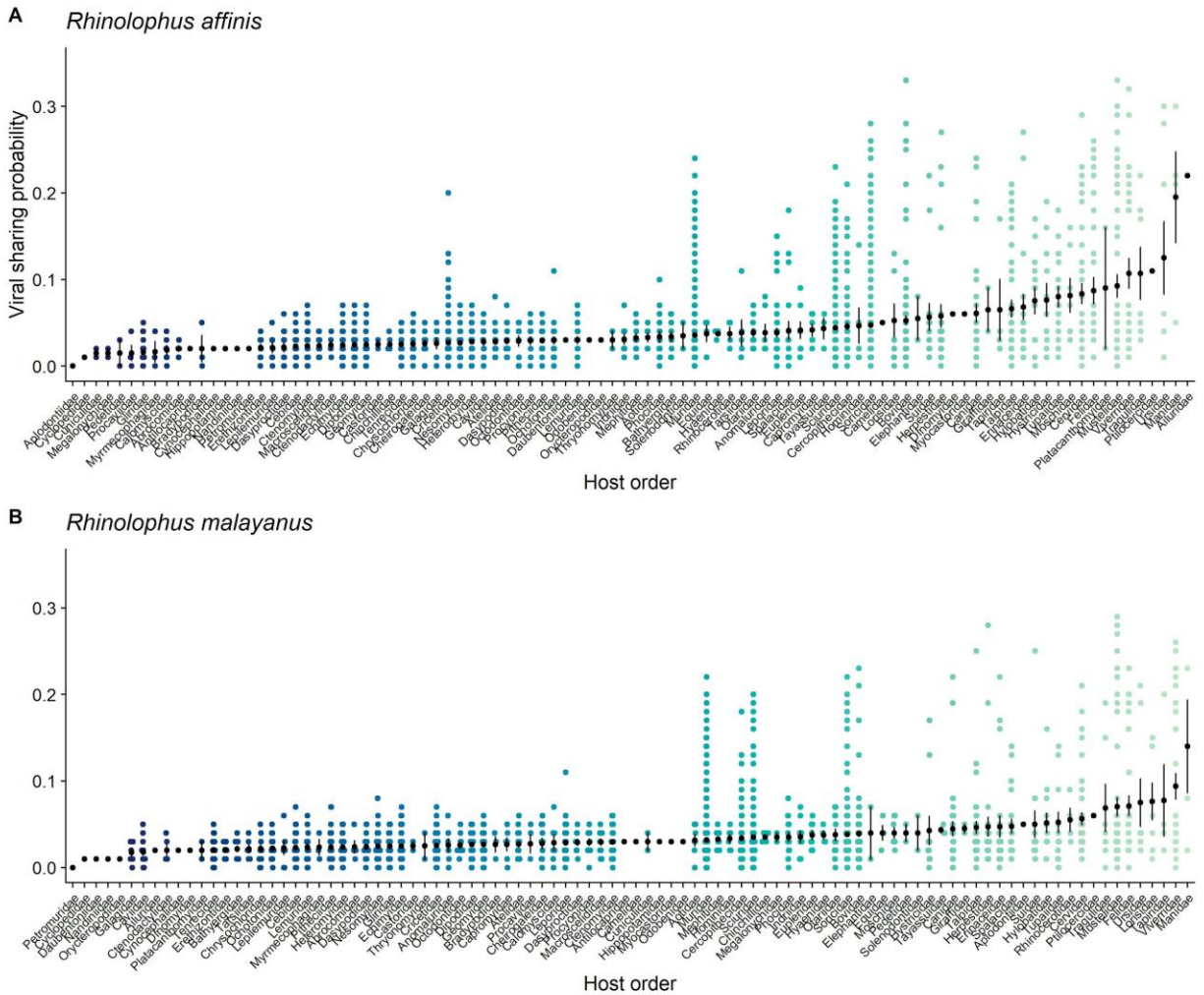
1013 **Extended Data Figure 4.** Predicted species-level sharing probabilities of A) *Rhinolophus affinis*  
1014 and B) *Rhinolophus malayanus*, calculated according to the phylogeographic viral sharing  
1015 model<sup>48</sup>. Each coloured point is a mammal species. Black points and error bars denote means  
1016 and standard errors for each order. Mammal orders are arranged according to their mean sharing  
1017 probability.  
1018



1019

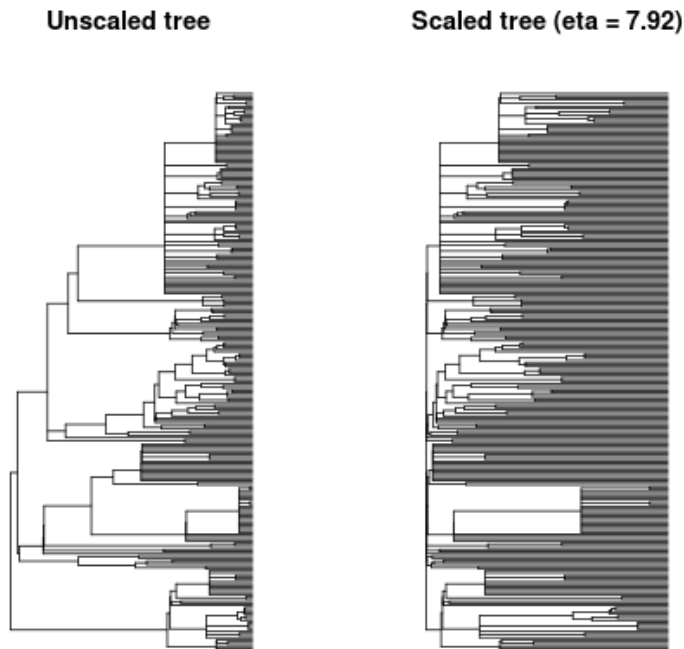


1020 **Extended Data Figure 5.** Predicted species-level sharing probabilities of A) *Rhinolophus affinis*  
1021 and B) *Rhinolophus malayanus*, calculated according to the phylogeographic viral sharing model<sup>4</sup>.  
1022 Each coloured point is a mammal species. Black points and error bars denote means and  
1023 standard errors for each order. Mammal orders are arranged according to their mean sharing  
1024 probability.  
1025



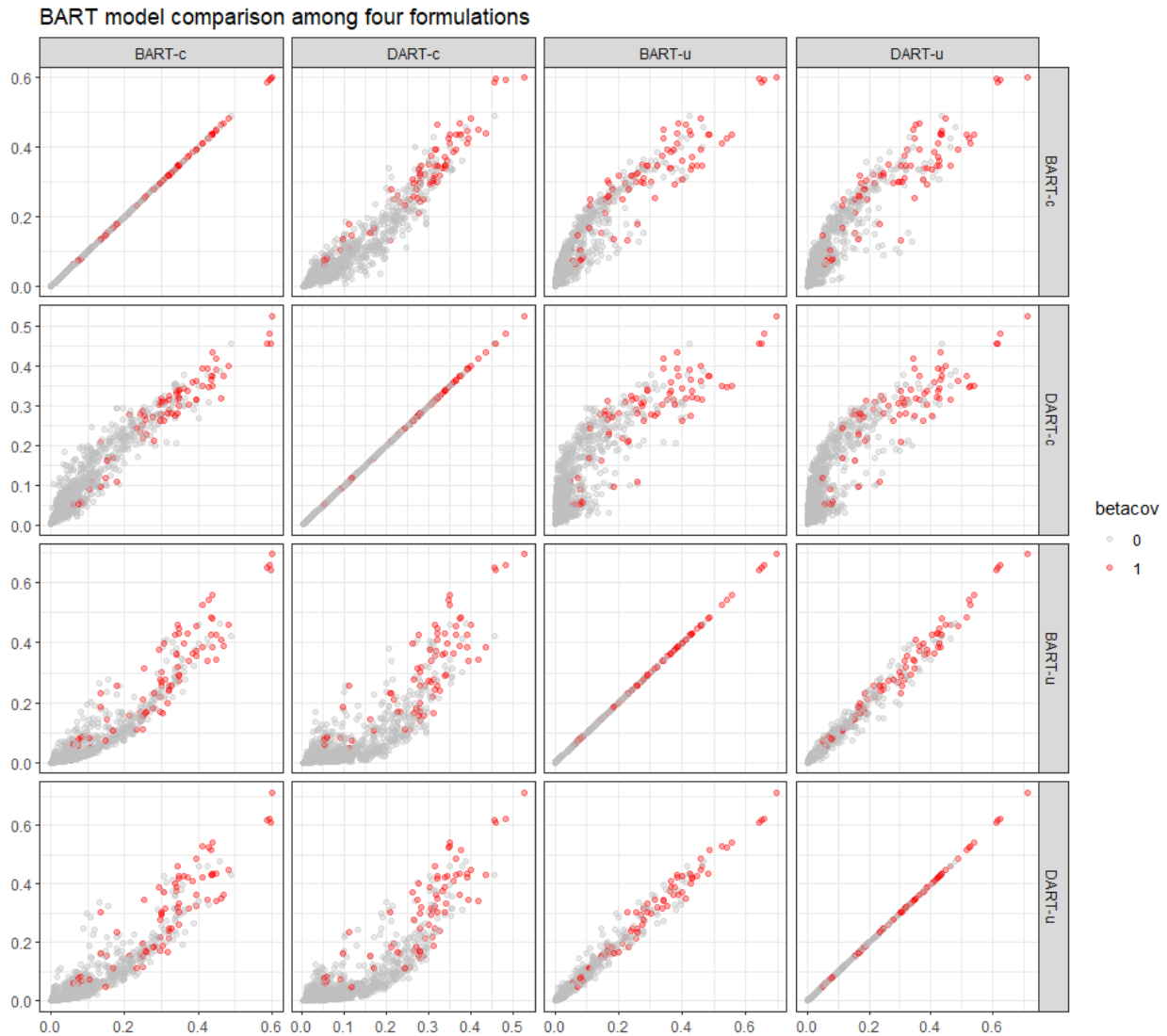
1026  
1027  
1028  
1029

1030 **Extended Data Figure 6.** To account for uncertainty in the phylogenetic distances among hosts,  
1031 the scaled-phylogeny model estimates a tree scaling parameter ( $\eta$ ) based on an early-burst  
1032 model of evolution. On the left is the unscaled bat phylogeny for the hosts in the bat-virus genera  
1033 network, and on the right is the same tree rescaled according to mean estimated scaling  
1034 parameter ( $\eta = 7.92$ ).  $\eta$  values above 1 indicate accelerating evolution, suggesting less  
1035 phylogenetic conservatism in host-virus associations among closely related taxa than would be  
1036 predicted by a Brownian-motion model on the unscaled tree.  
1037

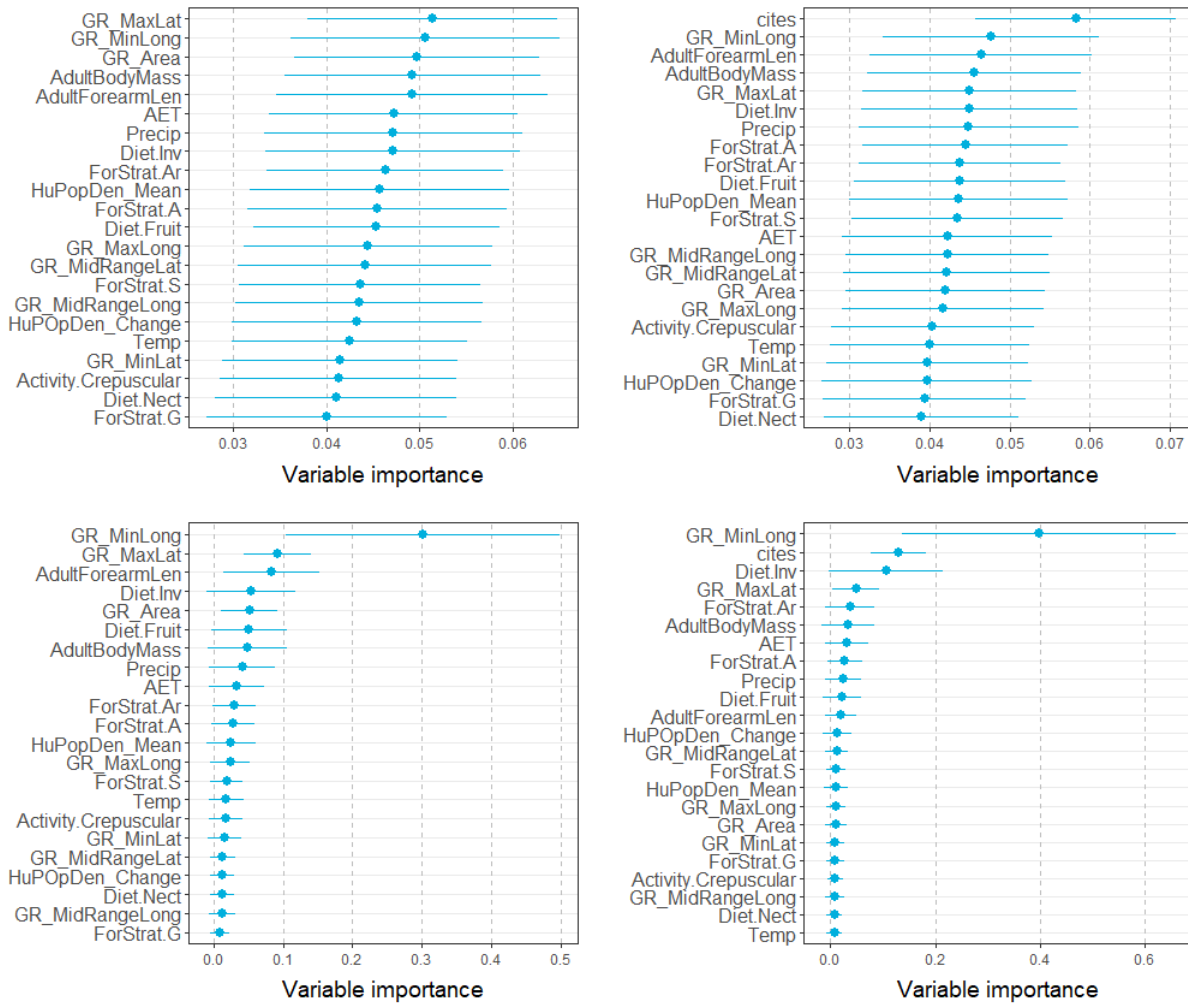


1038  
1039

1040 **Extended Data Figure 7.** Four formulations of Bayesian additive regression tree (BART) models  
1041 produce slightly different results, but largely agree. Two models use baseline BART, while two  
1042 models use a Dirichlet prior on variable importance (DART). Two are uncorrected for sampling  
1043 bias (u) while two are corrected using citation counts (c). In the final main-text model ensemble,  
1044 we present a DART model including correction for citation bias, which penalizes overfitting and  
1045 spurious patterns two ways and leads to predictions with a lower total correlation with the data,  
1046 but a still-high performance (AUC = 0.90).  
1047

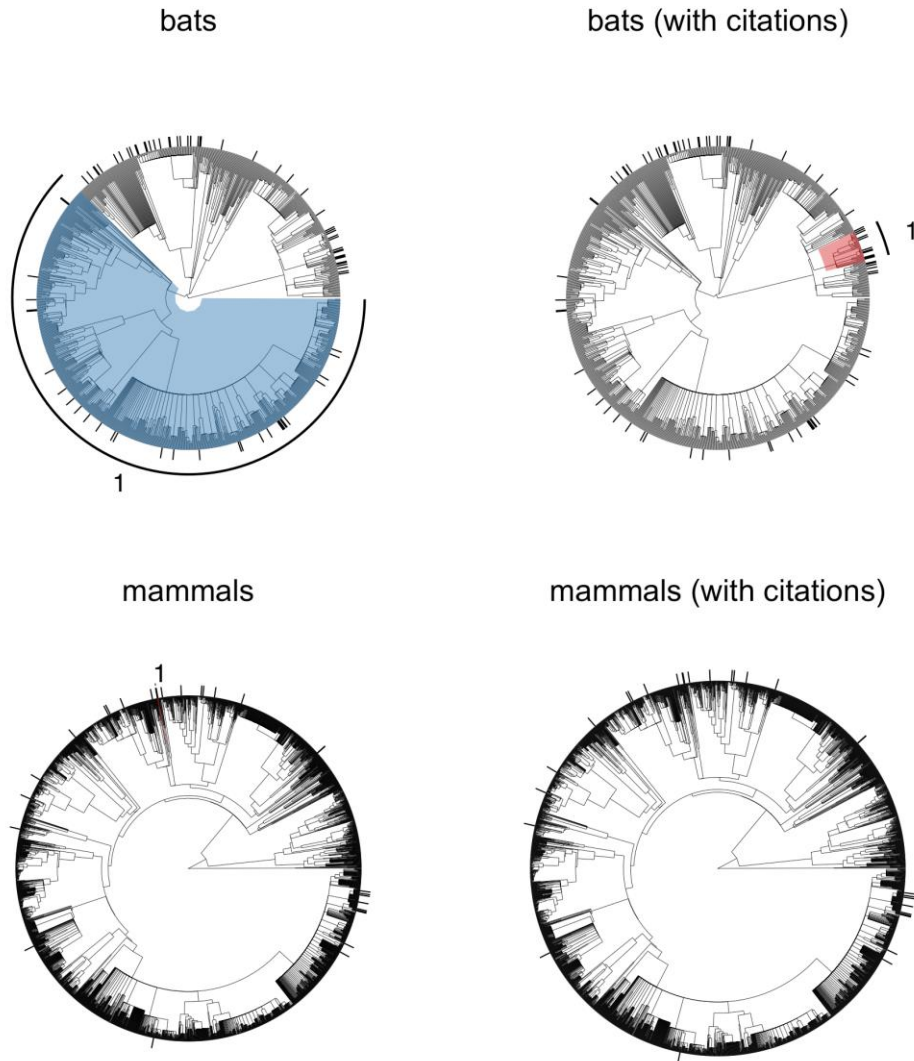


1049 **Extended Data Figure 8.** Partial dependence for the Bayesian additive regression tree models with  
 1050 uniform variable importance prior (top) versus Dirichlet prior (bottom), without (left) and with  
 1051 (right) correction for citations.



1052  
 1053  
 1054

1055 **Extended Data Figure 9.** Results of phylogenetic factorization applied to binomial  
1056 betacoronavirus data across bats (top) and other mammals (bottom), using raw data (left) and  
1057 after weighting by citation counts (right). Any significant clades (5% family-wise error rate) are  
1058 displayed in colored shading on the phylogeny. Bars indicate betacoronavirus detection, and  
1059 clades are colored by having more (red) or fewer (blue) positive species.  
1060



1061

1062  
1063  
1064

## Bibliography

- 1065 1. Anthony, S. J. *et al.* Global patterns in coronavirus diversity. *Virus Evol* **3**, vex012 (2017).
- 1066 2. Denison, M. R., Graham, R. L., Donaldson, E. F., Eckerle, L. D. & Baric, R. S. Coronaviruses: an  
1067 RNA proofreading machine regulates replication fidelity and diversity. *RNA Biol.* **8**, 270–279  
1068 (2011).
- 1069 3. Ren, W. *et al.* Full-length genome sequences of two SARS-like coronaviruses in horseshoe  
1070 bats and genetic variation analysis. *J. Gen. Virol.* **87**, 3355–3359 (2006).
- 1071 4. Li, W. *et al.* Bats are natural reservoirs of SARS-like coronaviruses. *Science* **310**, 676–679  
1072 (2005).
- 1073 5. Yang, X.-L. *et al.* Isolation and Characterization of a Novel Bat Coronavirus Closely Related  
1074 to the Direct Progenitor of Severe Acute Respiratory Syndrome Coronavirus. *J. Virol.* **90**,  
1075 3253–3256 (2015).
- 1076 6. Zhou, P. *et al.* A pneumonia outbreak associated with a new coronavirus of probable bat  
1077 origin. *Nature* **579**, 270–273 (2020).
- 1078 7. Guan, Y. *et al.* Isolation and characterization of viruses related to the SARS coronavirus  
1079 from animals in southern China. *Science* **302**, 276–278 (2003).
- 1080 8. Hu, B. *et al.* Discovery of a rich gene pool of bat SARS-related coronaviruses provides new  
1081 insights into the origin of SARS coronavirus. *PLoS Pathog.* **13**, e1006698 (2017).
- 1082 9. Memish, Z. A. *et al.* Middle East respiratory syndrome coronavirus in bats, Saudi Arabia.  
1083 *Emerg. Infect. Dis.* **19**, 1819–1823 (2013).
- 1084 10. Wang, Q. *et al.* Bat origins of MERS-CoV supported by bat coronavirus HKU4 usage of  
1085 human receptor CD26. *Cell Host Microbe* **16**, 328–337 (2014).
- 1086 11. Yang, Y. *et al.* Receptor usage and cell entry of bat coronavirus HKU4 provide insight into  
1087 bat-to-human transmission of MERS coronavirus. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 12516–  
1088 12521 (2014).
- 1089 12. Hu, B., Ge, X., Wang, L.-F. & Shi, Z. Bat origin of human coronaviruses. *Virology Journal* vol.  
1090 12 (2015).
- 1091 13. Anthony, S. J. *et al.* Further Evidence for Bats as the Evolutionary Source of Middle East  
1092 Respiratory Syndrome Coronavirus. *MBio* **8**, (2017).
- 1093 14. Anthony, S. J. *et al.* Coronaviruses in bats from Mexico. *J. Gen. Virol.* **94**, 1028–1038  
1094 (2013).
- 1095 15. Yang, L. *et al.* MERS-related betacoronavirus in *Vespertilio superans* bats, China. *Emerg.*  
1096 *Infect. Dis.* **20**, 1260–1262 (2014).
- 1097 16. Zhou, H. *et al.* A novel bat coronavirus reveals natural insertions at the S1/S2 cleavage site  
1098 of the Spike protein and a possible recombinant origin of HCoV-19.  
1099 doi:10.1101/2020.03.02.974139.
- 1100 17. Nielsen, R., Wang, H. & Pipes, L. Synonymous mutations and the molecular evolution of  
1101 SARS-Cov-2 origins. doi:10.1101/2020.04.20.052019.
- 1102 18. Lam, T. T.-Y. *et al.* Identifying SARS-CoV-2 related coronaviruses in Malayan pangolins.  
1103 *Nature* (2020) doi:10.1038/s41586-020-2169-0.
- 1104 19. Xiao, K. *et al.* Isolation of SARS-CoV-2-related coronavirus from Malayan pangolins. *Nature*  
1105 (2020) doi:10.1038/s41586-020-2313-x.
- 1106 20. Zhang, T., Wu, Q. & Zhang, Z. Probable Pangolin Origin of SARS-CoV-2 Associated with the  
1107 COVID-19 Outbreak. *Curr. Biol.* **30**, 1578 (2020).
- 1108 21. Andersen, K. G., Rambaut, A., Ian Lipkin, W., Holmes, E. C. & Garry, R. F. The proximal origin  
1109 of SARS-CoV-2. *Nature Medicine* vol. 26 450–452 (2020).



- 1110 22. Viana, M. *et al.* Assembling evidence for identifying reservoirs of infection. *Trends Ecol.*  
1111 *Evol.* **29**, 270–279 (2014).
- 1112 23. Plowright, R. K. *et al.* Prioritizing surveillance of Nipah virus in India. *PLoS Negl. Trop. Dis.*  
1113 **13**, e0007393 (2019).
- 1114 24. Becker, D. J., Crowley, D. E., Washburne, A. D. & Plowright, R. K. Temporal and spatial  
1115 limitations in global surveillance for bat filoviruses and henipaviruses. *Biol. Lett.* **15**,  
1116 20190423 (2019).
- 1117 25. Becker, D. J., Washburne, A. D., Faust, C. L., Mordecai, E. A. & Plowright, R. K. The problem  
1118 of scale in the prediction and management of pathogen spillover. *Philos. Trans. R. Soc.*  
1119 *Lond. B Biol. Sci.* **374**, 20190224 (2019).
- 1120 26. Becker, D. J. & Han, B. A. The macroecology and evolution of avian competence for *Borrelia*  
1121 *burgdorferi*. doi:10.1101/2020.04.15.040352.
- 1122 27. Han, B. A. *et al.* Undiscovered Bat Hosts of Filoviruses. *PLoS Negl. Trop. Dis.* **10**, e0004815  
1123 (2016).
- 1124 28. Han, B. A. *et al.* Confronting data sparsity to identify potential sources of Zika virus spillover  
1125 infection among primates. *Epidemics* **27**, 59–65 (2019).
- 1126 29. Washburne, A. D. *et al.* Taxonomic patterns in the zoonotic potential of mammalian viruses.  
1127 *PeerJ* **6**, e5979 (2018).
- 1128 30. Olival, K. J. *et al.* Host and viral traits predict zoonotic spillover from mammals. *Nature* **546**,  
1129 646–650 (2017).
- 1130 31. Fritz, S. A., Bininda-Emonds, O. R. P. & Purvis, A. Geographical variation in predictors of  
1131 mammalian extinction risk: big is bad, but only in the tropics. *Ecol. Lett.* **12**, 538–549  
1132 (2009).
- 1133 32. Jones, K. E. *et al.* PanTHERIA: a species-level database of life history, ecology, and  
1134 geography of extant and recently extinct mammals. *Ecology* vol. 90 2648–2648 (2009).
- 1135 33. Wilman, H. *et al.* EltonTraits 1.0: Species-level foraging attributes of the world’s birds and  
1136 mammals. *Ecology* vol. 95 2027–2027 (2014).
- 1137 34. Trifonova, N. *et al.* Spatio-temporal Bayesian network models with latent variables for  
1138 revealing trophic dynamics and functional networks in fisheries ecology. *Ecological*  
1139 *Informatics* vol. 30 142–158 (2015).
- 1140 35. Rohr, R. P., Scherer, H., Kehrli, P., Mazza, C. & Bersier, L.-F. Modeling food webs: exploring  
1141 unexplained structure using latent traits. *Am. Nat.* **176**, 170–177 (2010).
- 1142 36. Dallas, T., Park, A. W. & Drake, J. M. Predicting cryptic links in host-parasite networks. *PLoS*  
1143 *Comput. Biol.* **13**, e1005557 (2017).
- 1144 37. Han, B. A., Schmidt, J. P., Bowden, S. E. & Drake, J. M. Rodent reservoirs of future zoonotic  
1145 diseases. *Proceedings of the National Academy of Sciences* vol. 112 7039–7044 (2015).
- 1146 38. Brandão, P. E. *et al.* A coronavirus detected in the vampire bat *Desmodus rotundus*. *Braz. J.*  
1147 *Infect. Dis.* **12**, 466–468 (2008).
- 1148 39. Corman, V. M. *et al.* Highly diversified coronaviruses in neotropical bats. *J. Gen. Virol.* **94**,  
1149 1984–1994 (2013).
- 1150 40. Moreira-Soto, A. *et al.* Neotropical bats from Costa Rica harbour diverse coronaviruses.  
1151 *Zoonoses Public Health* **62**, 501–505 (2015).
- 1152 41. Wang, L. *et al.* Discovery and genetic analysis of novel coronaviruses in least horseshoe  
1153 bats in southwestern China. *Emerg. Microbes Infect.* **6**, e14 (2017).
- 1154 42. Lin, X.-D. *et al.* Extensive diversity of coronaviruses in bats from China. *Virology* vol. 507 1–  
1155 10 (2017).
- 1156 43. Wacharapluesadee, S. *et al.* Diversity of coronavirus in bats from Eastern Thailand. *Viol. J.*  
1157 **12**, 57 (2015).
- 1158 44. Guy, C., Ratcliffe, J. M. & Mideo, N. The influence of bat ecology on viral diversity and

- 1159 reservoir status. *Ecol. Evol.* **2008**, 209 (2020).
- 1160 45. Washburne, A. D. et al. Phylofactorization: a graph partitioning algorithm to identify  
1161 phylogenetic scales of ecological data. *Ecol. Monogr.* **89**, e01353 (2019).
- 1162 46. Almeida, F. C., Simmons, N. B. & Giannini, N. P. A Species-Level Phylogeny of Old World  
1163 Fruit Bats with a New Higher-Level Classification of the Family Pteropodidae. *American*  
1164 *Museum Novitates* vol. 2020 1 (2020).
- 1165 47. Crowley, D., Becker, D., Washburne, A. & Plowright, R. Identifying Suspect Bat Reservoirs of  
1166 Emerging Infections. *Vaccines* vol. 8 228 (2020).
- 1167 48. Albery, G. F., Eskew, E. A., Ross, N. & Olival, K. J. Predicting the global mammalian viral  
1168 sharing network using phylogeography. *Nat. Commun.* **11**, 2260 (2020).
- 1169 49. Wang, M. et al. SARS-CoV Infection in a Restaurant from Palm Civet. *Emerging Infectious*  
1170 *Diseases* vol. 11 1860–1865 (2005).
- 1171 50. Song, H.-D. et al. Cross-host evolution of severe acute respiratory syndrome coronavirus in  
1172 palm civet and human. *Proc. Natl. Acad. Sci. U. S. A.* **102**, 2430–2435 (2005).
- 1173 51. Oreshkova, N. et al. SARS-CoV2 infection in farmed mink, Netherlands, April 2020.  
1174 doi:10.1101/2020.05.18.101493.
- 1175 52. Damas, J., Hughes, G. M., Keough, K. C. & Painter, C. A. Broad Host Range of SARS-CoV-2  
1176 Predicted by Comparative and Structural Analysis of ACE2 in Vertebrates. *bioRxiv* (2020).
- 1177 53. Shi, J. et al. Susceptibility of ferrets, cats, dogs, and other domesticated animals to SARS-  
1178 coronavirus 2. *Science* (2020) doi:10.1126/science.abb7015.
- 1179 54. Yang, X.-L. et al. Genetically Diverse Filoviruses in Rousettus and Eonycteris spp. Bats,  
1180 China, 2009 and 2015. *Emerg. Infect. Dis.* **23**, 482–486 (2017).
- 1181 55. Seifert, S. N. et al. Rousettus aegyptiacus Bats Do Not Support Productive Nipah Virus  
1182 Replication. *The Journal of Infectious Diseases* vol. 221 S407–S413 (2020).
- 1183 56. Wacharapluesadee, S. et al. Longitudinal study of age-specific pattern of coronavirus  
1184 infection in Lyle's flying fox (*Pteropus lylei*) in Thailand. *Viol. J.* **15**, 38 (2018).
- 1185 57. Yang, L. et al. Novel SARS-like betacoronaviruses in bats, China, 2011. *Emerg. Infect. Dis.*  
1186 **19**, 989–991 (2013).
- 1187 58. Geldenhuys, M. et al. A metagenomic viral discovery approach identifies potential zoonotic  
1188 and novel mammalian viruses in Neoromicia bats within South Africa. *PLoS One* **13**,  
1189 e0194527 (2018).
- 1190 59. Memish, Z. A. et al. Middle East respiratory syndrome coronavirus in bats, Saudi Arabia.  
1191 *Emerg. Infect. Dis.* **19**, 1819–1823 (2013).
- 1192 60. Luo, Y. et al. Longitudinal Surveillance of Betacoronaviruses in Fruit Bats in Yunnan  
1193 Province, China During 2009-2016. *Viol. Sin.* **33**, 87–95 (2018).
- 1194 61. Peel, A. J. et al. Synchronous shedding of multiple bat paramyxoviruses coincides with  
1195 peak periods of Hendra virus spillover. *Emerg. Microbes Infect.* **8**, 1314–1323 (2019).
- 1196 62. de Souza Cortez J. L. Dunnum A. W. Ferguson F. A. Anwarali Khan D. L. Paul D. M. Reeder  
1197 N. B. Simmons B. M. Thiers C. W. Thompson N S. Upham M. P. M. Vanhove P. W. Webala  
1198 M. Weksler R. Yanagihara P. S. Soltis., C. J. A. S. A. B. A. J. B. C. A. C. B. M. B. Integrating  
1199 biodiversity infrastructure into pathogen discovery and mitigation of epidemic infectious  
1200 diseases. *Bioscience* (2020) doi:biaa064.
- 1201 63. Kingston, T. et al. Networking networks for global bat conservation. in *Bats in the*  
1202 *Anthropocene: Conservation of Bats in a Changing World* 539–569 (Springer, Cham, 2016).
- 1203 64. Phelps, K. L. et al. Bat Research Networks and Viral Surveillance: Gaps and Opportunities in  
1204 Western Asia. *Viruses* **11**, (2019).
- 1205 65. Teeling, E. C. et al. Bat Biology, Genomes, and the Bat1K Project: To Generate  
1206 Chromosome-Level Genomes for All Living Bat Species. *Annu Rev Anim Biosci* **6**, 23–46  
1207 (2018).

- 1208 66. Mandl, J. N., Schneider, C., Schneider, D. S. & Baker, M. L. Going to Bat(s) for Studies of  
1209 Disease Tolerance. *Front. Immunol.* **9**, 2112 (2018).
- 1210 67. Gervasi, S. S., Civitello, D. J., Kilvitis, H. J. & Martin, L. B. The context of host competence: a  
1211 role for plasticity in host–parasite dynamics. *Trends Parasitol.* **31**, 419–425 (2015).
- 1212 68. Martin, L. B., Burgan, S. C., Adelman, J. S. & Gervasi, S. S. Host Competence: An Organismal  
1213 Trait to Integrate Immunology and Epidemiology. *Integr. Comp. Biol.* **56**, 1225–1237 (2016).
- 1214 69. Callaway, E. & Cyranoski, D. Why snakes probably aren't spreading the new China virus.  
1215 *Nature* (2020) doi:10.1038/d41586-020-00180-8.
- 1216 70. Gong, Y., Wen, G., Jiang, J. & Feng, X. Complete title: Codon bias analysis may be  
1217 insufficient for identifying host(s) of a novel virus. *J. Med. Virol.* (2020)  
1218 doi:10.1002/jmv.25977.
- 1219 71. Zhao, H. COVID-19 drives new threat to bats in China. *Science* **367**, 1436 (2020).
- 1220 72. Fenton, M. B. *et al.* Knowledge gaps about rabies transmission from vampire bats to  
1221 humans. *Nature Ecology & Evolution* **4**, 517–518 (2020).
- 1222 73. López-Baucells, A., Rocha, R. & Fernández-Llamazares, Á. When bats go viral: negative  
1223 framings in virological research imperil bat conservation. *Mamm. Rev.* **48**, 62–66 (2018).
- 1224 74. O'Shea, T. J., Cryan, P. M., Hayman, D. T. S., Plowright, R. K. & Streicker, D. G. Multiple  
1225 mortality events in bats: a global review. *Mamm. Rev.* **46**, 175–190 (2016).
- 1226 75. MB Fenton, S Mubareka, SM Tsang, NB Simmons, DJ Becker. COVID-19 and threats to bats.  
1227 *FACETS* in press (2020).
- 1228 76. Aguiar, L. M. S., Brito, D. & Machado, R. B. Do current vampire bat (*Desmodus rotundus*)  
1229 population control practices pose a threat to Dekeyser's nectar bat's (*Lonchophylla*  
1230 *dekeyseri*) long-term persistence in the Cerrado? *Acta Chiropt.* **12**, 275–282 (2010).
- 1231 77. Streicker, D. G. *et al.* Ecological and anthropogenic drivers of rabies exposure in vampire  
1232 bats: implications for transmission and control. *Proc. Biol. Sci.* **279**, 3384–3392 (2012).
- 1233 78. Blackwood, J. C., Streicker, D. G., Altizer, S. & Rohani, P. Resolving the roles of immunity,  
1234 pathogenesis, and immigration for rabies persistence in vampire bats. *Proc. Natl. Acad. Sci.*  
1235 *U. S. A.* **110**, 20837–20842 (2013).
- 1236 79. Frick, W. F. *et al.* An emerging disease causes regional population collapse of a common  
1237 North American bat species. *Science* **329**, 679–682 (2010).
- 1238 80. Frick, W. F. *et al.* Disease alters macroecological patterns of North American bats. *Glob.*  
1239 *Ecol. Biogeogr.* **24**, 741–749 (2015).
- 1240 81. Sabir, J. S. M. *et al.* Co-circulation of three camel coronavirus species and recombination of  
1241 MERS-CoVs in Saudi Arabia. *Science* **351**, 81–84 (2016).
- 1242 82. Guth, S., Visher, E., Boots, M. & Brook, C. E. Host phylogenetic distance drives trends in virus  
1243 virulence and transmissibility across the animal–human interface. *Philos. Trans. R. Soc.*  
1244 *Lond. B Biol. Sci.* **374**, 20190296 (2019).
- 1245 83. Redondo, R. A. F., Brina, L. P. S., Silva, R. F., Ditchfield, A. D. & Santos, F. R. Molecular  
1246 systematics of the genus *Artibeus* (Chiroptera: Phyllostomidae). *Mol. Phylogenet. Evol.* **49**,  
1247 44–58 (2008).
- 1248 84. Bouchard, S. *Chaerephon pumilus*. *Mammalian Species* 1–6 (1998).
- 1249 85. Hooper, S. R., Van Den Bussche, R. A. & Horáček, I. Generic Status of the American  
1250 Pipistrelles (Vespertilionidae) with Description of a New Genus. *J. Mammal.* **87**, 981–992  
1251 (2006).
- 1252 86. Desjardins-Proulx, P., Laigle, I., Poisot, T. & Gravel, D. Ecological interactions and the Netflix  
1253 problem. *PeerJ* **5**, e3644 (2017).
- 1254 87. Stock, M., Poisot, T., Waegeman, W. & De Baets, B. Linear filtering reveals false negatives in  
1255 species interaction data. *Sci. Rep.* **7**, 45908 (2017).
- 1256 88. Drake, J. M. & Richards, R. L. Estimating environmental suitability. *Ecosphere* vol. 9 e02373

- 1257 (2018).
- 1258 89. Dallas, T. A., Carlson, C. J. & Poisot, T. Testing predictability of disease outbreaks with a  
1259 simple model of pathogen biogeography. *R Soc Open Sci* **6**, 190883 (2019).
- 1260 90. Elmasri, M., Farrell, M. J., Jonathan Davies, T. & Stephens, D. A. A hierarchical Bayesian  
1261 model for predicting ecological interactions using scaled evolutionary relationships. *The*  
1262 *Annals of Applied Statistics* vol. 14 221–240 (2020).
- 1263 91. Cadotte, M. W. *et al.* Phylogenetic diversity metrics for ecological communities: integrating  
1264 species richness, abundance and evolutionary history. *Ecol. Lett.* **13**, 96–105 (2010).
- 1265 92. Park, A. W. *et al.* Characterizing the phylogenetic specialism–generalism spectrum of  
1266 mammal parasites. *Proceedings of the Royal Society B: Biological Sciences* vol. 285  
1267 20172613 (2018).
- 1268 93. Harvey, P. H. & Pagel, M. D. *The comparative method in evolutionary biology*. (Oxford  
1269 University Press, USA, 1998).
- 1270 94. Harmon, L. J. *et al.* Early bursts of body size and shape evolution are rare in comparative  
1271 data. *Evolution* **64**, 2385–2396 (2010).
- 1272 95. Mollentze, N. & Streicker, D. G. Viral zoonotic risk is homogenous among taxonomic orders  
1273 of mammalian and avian reservoir hosts. *Proc. Natl. Acad. Sci. U. S. A.* **117**, 9423–9430  
1274 (2020).
- 1275 96. Schmidt, J. P. *et al.* Ecological indicators of mammal exposure to Ebolavirus. *Philos. Trans.*  
1276 *R. Soc. Lond. B Biol. Sci.* **374**, 20180337 (2019).
- 1277 97. Pandit, P. S. *et al.* Predicting wildlife reservoirs and global vulnerability to zoonotic  
1278 Flaviviruses. *Nat. Commun.* **9**, 5425 (2018).
- 1279 98. Evans, M. V., Dallas, T. A., Han, B. A., Murdock, C. C. & Drake, J. M. Data-driven identification  
1280 of potential Zika virus vectors. *Elife* **6**, (2017).
- 1281 99. Yang, L. H. & Han, B. A. Data-driven predictions and novel hypotheses about zoonotic tick  
1282 vectors from the genus Ixodes. *BMC Ecol.* **18**, 7 (2018).
- 1283 100. Elith, J., Leathwick, J. R. & Hastie, T. A working guide to boosted regression trees. *J. Anim.*  
1284 *Ecol.* **77**, 802–813 (2008).
- 1285 101. Carlson, C. J. *embarcadero*: Species distribution modelling with Bayesian additive  
1286 regression trees in R. doi:10.1101/774604.
- 1287 102. Chipman, H. A., George, E. I. & McCulloch, R. E. BART: Bayesian additive regression trees.  
1288 *The Annals of Applied Statistics* vol. 4 266–298 (2010).
- 1289 103. Chen, W. *et al.* The illegal exploitation of hog badgers (*Arctonyx collaris*) in China: genetic  
1290 evidence exposes regional population impacts. *Conservation Genetics Resources* vol. 7  
1291 697–704 (2015).
- 1292 104. Foley, N. M., Goodman, S. M., Whelan, C. V., Puechmaille, S. J. & Teeling, E. Towards  
1293 navigating the Minotaur’s labyrinth: cryptic diversity and taxonomic revision within the  
1294 speciose genus *Hipposideros* (Hipposideridae). *Acta Chiropt.* **19**, 1–18 (2017).

1295

1296