

33 **Abstract.**

34

35 Despite massive investment in research on reservoirs of emerging pathogens, it remains difficult to
36 rapidly identify the wildlife origins of novel zoonotic viruses. Viral surveillance is costly but rarely
37 optimized using model-guided prioritization strategies, and predictions from a single model may be
38 highly uncertain. Here, we generate an ensemble of eight network- and trait-based statistical models that
39 predict mammal-virus associations, and we use model predictions to develop a set of priority
40 recommendations for sampling potential bat reservoirs and intermediate hosts for SARS-CoV-2 and
41 related betacoronaviruses. We find over 200 bat species globally could be undetected hosts of
42 betacoronaviruses. Although over a dozen species of Asian horseshoe bats (*Rhinolophus* spp.) are known
43 to harbor SARS-like coronaviruses, we find at least two thirds of betacoronavirus reservoirs in this bat
44 genus might still be undetected. Although identification of other probable mammal reservoirs is likely
45 beyond existing predictive capacity, some of our findings are surprisingly plausible; for example, several
46 civet and pangolin species were highlighted as high-priority species for viral sampling. Our results should
47 not be over-interpreted as novel information about the plausibility or likelihood of SARS-CoV-2's
48 ultimate origin, but rather these predictions could help guide sampling for novel potentially zoonotic
49 viruses; immunological research to characterize key receptors (e.g., ACE2) and identify mechanisms of
50 viral tolerance; and experimental infections to quantify competence of suspected host species.

51 Introduction

52

53 Coronaviruses are a diverse family of positive-sense, single-stranded RNA viruses, found widely in
54 mammals and birds¹. They have a broad host range, a high mutation rate, and the largest genomes of any
55 RNA viruses, but they have also evolved mechanisms for RNA proofreading and repair, which help to
56 mitigate the deleterious effects of a high recombination rate acting over a large genome². Consequently,
57 coronaviruses fit the profile of viruses with high zoonotic potential. There are seven human coronaviruses
58 (two in the genus *Alphacoronavirus* and five in *Betacoronavirus*), of which three are highly pathogenic in
59 humans: SARS-CoV, SARS-CoV-2, and MERS-CoV. These three are zoonotic and widely agreed to
60 have evolutionary origins in bats³⁻⁶.

61

62 Our collective experience with both SARS-CoV and MERS-CoV illustrate the difficulty of tracing
63 specific animal hosts of emerging coronaviruses. During the 2002–2003 SARS epidemic, SARS-CoV
64 was traced to the masked palm civet (*Paguma larvata*)⁷, but the ultimate origin remained unknown for
65 several years. Horseshoe bats (family Rhinolophidae: *Rhinolophus*) were implicated as reservoir hosts in
66 2005, but their SARS-like viruses were not identical to circulating human strains⁴. Stronger evidence
67 from 2017 placed the most likely evolutionary origin of SARS-CoV in *Rhinolophus ferrumequinum* or
68 potentially *R. sinicus*⁸. Presently, there is even less certainty in the origins of MERS-CoV, although
69 spillover to humans occurs relatively often through contact with dromedary camels (*Camelus*
70 *dromedarius*). A virus with 100% nucleotide identity in a ~200 base pair region of the polymerase gene
71 was detected in *Taphozous* bats (family Emballonuridae) in Saudi Arabia⁹; however, based on spike gene
72 similarity, other sources treat HKU4 virus from *Tylonycteris* bats (family Vespertilionidae) in China as
73 the closest-related bat virus^{10,11}. Several bat coronaviruses have shown close relation to MERS-CoV, with
74 a surprisingly broad geographic distribution from Mexico to China^{12,13,14,15}.

75

76 Coronavirus disease 2019 (COVID-19) is caused by severe acute respiratory syndrome coronavirus-2
77 (SARS-CoV-2), a novel virus with presumed evolutionary origins in bats. Although the earliest cases
78 were linked to a wildlife market, contact tracing was limited, and there has been no definitive
79 identification of the wildlife contact that resulted in spillover nor a true “index case.” Two bat viruses are
80 closely related to SARS-CoV-2: RaTG13 bat CoV from *Rhinolophus affinis* (96% identical overall), and
81 RmYN02 bat CoV from *Rhinolophus malayanus* (97% identical in one gene but only 61% in the
82 receptor-binding domain and with less overall similarity)^{6,16}. The divergence time between these bat
83 viruses and human SARS-CoV-2 has been estimated as 40-50 years¹⁷, suggesting that the main host(s)
84 involved in spillover remain unknown. Evidence of viral recombination in pangolins has been proposed
85 but is unresolved¹⁷. SARS-like betacoronaviruses have been recently isolated from Malayan pangolins
86 (*Manis javanica*) traded in wildlife markets^{18,19}, and these viruses have a very high amino acid identity to
87 SARS-CoV-2, but only show a ~90% nucleotide identity with SARS-CoV-2 or Bat-CoV RaTG13²⁰. None
88 of these host species are universally accepted as the origin of SARS-CoV-2 or a progenitor virus, and a
89 “better fit” wildlife reservoir could likely still be identified. However, substantial gaps in betacoronavirus
90 sampling across wildlife limit actionable inference about plausible reservoirs and intermediate hosts for
91 SARS-CoV-2²¹.

92

93 Identifying likely reservoirs of zoonotic pathogens is challenging²². Sampling wildlife for the presence of
94 active or previous infection (i.e., seropositivity) represents the first stage of a pipeline for proper inference

95 of host species²³, but sampling is often limited in phylogenetic, temporal, and spatial scale by logistical
96 constraints²⁴. Given such restrictions, modeling efforts can play a critical role in helping to prioritize
97 pathogen surveillance by narrowing the set of plausible sampling targets²⁵. For example, machine learning
98 approaches have generated candidate lists of likely, but unsampled, primate reservoirs for Zika virus, bat
99 reservoirs for filoviruses, and avian reservoirs for *Borrelia burgdorferi*^{26–28}. In some contexts, models
100 may be more useful for identifying which host or pathogen groups are *unlikely* to have zoonotic
101 potential²⁹. However, these approaches are generally applied individually to generate predictions.
102 Implementation of multiple modeling approaches collaboratively and simultaneously could reduce
103 redundancy and apparent disagreement at the earliest stages of pathogen tracing and help advance
104 modeling work by addressing inter-model reliability, predictive accuracy, and the broader utility (or
105 inefficacy) of such models in zoonosis research.

106
107 Because SARS-like coronaviruses (subgenus *Sarbecovirus*) are only characterized from a small number
108 of bat species in publicly available data, current modeling methods are poorly tailored to exactly infer
109 their potential reservoir hosts. In this study, we instead conduct two predictive efforts that may help guide
110 the inevitable search for known and future zoonotic coronaviruses in wildlife: (1) broadly identifying bats
111 and other mammals that may host any *Betacoronavirus* and (2) specifically identifying species with a
112 high viral sharing probability with the two *Rhinolophus* species carrying the closest known wildlife
113 relatives of SARS-CoV-2. To do this, we developed a standardized dataset of mammal-virus associations
114 by integrating a previously published mammal-virus dataset³⁰ with a targeted scrape of all GenBank
115 accessions for Coronaviridae and their associated hosts. Our final dataset spanned 710 host species and
116 359 virus genera, including 107 mammal hosts of betacoronaviruses as well as hundreds of other (non-
117 coronavirus) association records. We integrated our host-virus data with a mammal phylogenetic
118 supertree³¹ and over 60 ecological traits of bat species^{27,32,33}.

119
120 We used these standardized data to generate an ensemble of predictive models. We drew on two popular
121 approaches to identify candidate reservoirs and intermediate hosts of betacoronaviruses. *Network-based*
122 *methods* estimate a full set of “true” unobserved host-virus interactions based on a recorded network of
123 associations (here, pairs of host species and associated viral genera). These methods are increasingly
124 popular as a way to identify latent processes structuring ecological networks^{34–36}, but they are often
125 confounded by sampling bias and can only make predictions for species within the observed network (i.e.,
126 those that have available virus data; in-sample prediction). In contrast, *trait-based methods* use observed
127 relationships concerning host traits to identify species that fit the morphological, ecological, and/or
128 phylogenetic profile of known host species of a given pathogen and rank the suitability of unknown hosts
129 based on these trait profiles^{28,37}. These methods may be more likely to recapitulate patterns in observed
130 host-pathogen association data (e.g., geographic biases in sampling, phylogenetic similarity in host
131 morphology), but they more easily correct for sampling bias and can predict host species without known
132 viral associations (out-of-sample prediction).

133
134 In total, we implemented eight different models of host-virus associations, including four network-based
135 approaches, three trait-based approaches, and one hybrid approach using both network and trait
136 information. These efforts generated eight ranked lists of suspected bat hosts of betacoronaviruses and
137 five ranked lists for other mammals. Each ranked list was scaled proportionally and consolidated in an

138 ensemble of recommendations for betacoronavirus sampling and broader eco-evolutionary research
139 (Supplemental Figure 1).

140

141 **Results**

142

143 *Predicted bat hosts of betacoronaviruses*

144

145 Predictions of bat hosts of betacoronavirus derived from network- and trait-based approaches displayed
146 strong inter-model agreement within-group, but less with each other (Figure 1A,B). In-sample, we
147 identified bat species across a range of genera as having the highest predicted probabilities of hosting
148 betacoronaviruses, distributed in distinct families in both the Old World (e.g., Hipposideridae, several
149 subfamilies in the Vespertilionidae) and the New World (e.g., *Artibeus jamaicensis* from the
150 Phyllostomidae; Figure 1C). Out-of-sample, our multi-model ensemble more conservatively limited
151 predictions to primarily Old World families such as Rhinolophidae and Pteropodidae (Figure 1D). Of the
152 1,037 bat hosts not currently known to host betacoronaviruses, our models identified between 1 and 720
153 potential hosts based on a 10% omission threshold (90% sensitivity). Applying this same threshold to our
154 ensemble predictions, we identified 239 bat species that are likely undetected hosts of betacoronaviruses.
155 These include more half of bat species in the genus *Rhinolophus* not currently known to be
156 betacoronavirus hosts (35 of 61), compared to 16 known hosts in this genus. Given known roles of
157 rhinolophids as hosts of SARS-like coronaviruses, our results suggest that SARS-like coronavirus
158 diversity could be undescribed for around two-thirds of the potential reservoir bat species.

159

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

166

167 Because only trait-based models were capable of out-of-sample prediction, the differences in geographic
168 patterns of our predictions likely reflect distinctions between the network- and trait-based modeling
169 approaches, which we suggest should be considered qualitatively different lines of evidence. Network
170 approaches proportionally upweight species with high observed viral diversity, recapitulating sampling
171 biases largely unrelated to coronaviruses (e.g., frequent screening for rabies lyssaviruses in vampire bats,
172 which have been sampled in a comparatively limited capacity for coronaviruses^{14,38-40}). Highly ranked
173 species may also have been previously sampled without evidence of betacoronavirus presence; for
174 example, *Rhinolophus luctus* and *Macroglossus sobrinus* from China and Thailand, respectively, tested
175 negative for betacoronaviruses, but detection probability was limited by small sample sizes⁴¹⁻⁴³. In
176 contrast, trait-based approaches are constrained by their reliance on phylogeny and ecological traits, and
177 the use of geographic covariates made models more likely to recapitulate existing spatial patterns of
178 betacoronavirus detection (i.e., clustering in southeast Asia). However, their out-of-sample predictions
179 are, by definition, inclusive of unsampled hosts⁴⁴, which potentially offer greater return on viral
180 discovery investment.

181

182 Multi-model ensemble predictions also clustered taxonomically along parallel lines. Applying a graph
183 partitioning algorithm (phylogenetic factorization) to the bat phylogeny⁴⁵, we found that in-sample
184 predictions were on average lowest for the Yangochiroptera (Figure 3). This makes intuitive sense,
185 because this clade does not include the groups known to harbor the majority of betacoronaviruses
186 detected in bats (e.g., *Rhinolophus*, Hipposideridae). Out-of-sample predictions were lower in the New
187 World superfamily Noctilionoidea, the emballonurids, and the *Lasiurus* genus, whereas the *Rhinolophus*
188 genus and the Old World fruit bats (Pteropodidae) both had higher mean probabilities of betacoronavirus
189 hosting (Supplemental Table 1).

190
191 These clade-specific patterns of predicted probabilities across extant bats could be particularly applicable
192 for guiding future surveillance. On the one hand, betacoronavirus sampling in southeast Asian bat taxa
193 (especially the genus *Rhinolophus*) may have a high success of viral detection but may not improve
194 existing bat sampling gaps⁴⁶. On the other hand, discovery of novel betacoronaviruses in Neotropical bats
195 or Old World fruit bats could significantly revise our understanding of the bat-virus association network.
196 Such discoveries would be particularly important for global health security, given the surprising
197 identification of a MERS-like virus in Mexican bats¹⁴ and the likelihood that post-COVID pandemic
198 preparedness efforts will focus disproportionately on Asia despite the near-global presence of bat
199 betacoronaviruses.

200

201 ***Mammal-wide predictions***

202

203 Although our ensemble model of potential bat betacoronavirus reservoirs generated strong and actionable
204 predictions, our mammal-wide predictions were largely uninformative. In particular, minimal inter-model
205 agreement (Supplemental Figure 2) indicated a lack of consistent, biologically meaningful findings.
206 Major effects of sampling bias were apparent from the top-ranked species, which were primarily domestic
207 animals or well-studied mesocarnivores (ED Figure 2B). Phylogenetic factorization mostly failed to find
208 specific patterns in prediction (Supplemental Table 2): in-sample, mean predictions primarily confirmed
209 betacoronavirus detection in the remaining Laurasiatheria (e.g., ungulates, carnivores, pangolins,
210 hedgehogs, shrews), although nested clades of marine mammals (i.e., cetaceans) were less likely to harbor
211 these viruses. Although cetaceans are predicted to be susceptible to SARS-CoV-2 based on ACE2
212 similarity with humans⁴⁷, this result is expected given betacoronavirus epidemiology and their
213 predominance in terrestrial mammals. Our mammal predictions thus reflect a combination of detection
214 bias and poor performance of network methods on limited data that likely signals the limits of existing
215 predictive capacity. Our dataset contained only 30 non-bat betacoronavirus hosts, many of which were
216 identified during sampling efforts following the first SARS outbreak⁷. Although the laurasiatherians
217 include more potential intermediate hosts than other mammals, likely driven by inclusion of bats in this
218 group, the high species diversity of this clade restricts insights for sampling prioritization, experimental
219 work, or spillover risk management.

220

221 ***SARS-CoV-2 sampling priorities***

222

223 Given the unresolved origins of SARS-CoV-2 and significant motivation to identify other SARS-like
224 coronaviruses and their reservoir hosts for pandemic preparedness²¹, we further explored our only model
225 that could generate out-of-sample predictions for all mammals⁴⁸. This model uses geographic distributions

226 and phylogenetic relatedness to estimate viral sharing probability. Where one or more (potential) hosts are
227 known, these sharing patterns can be interpreted to identify probable reservoir hosts⁴⁸. Because
228 *Rhinolophus affinis* and *R. malayanus* host viruses that are closely related to SARS-CoV-2^{6,16}, we used
229 their predicted sharing patterns to identify possible reservoirs of sarbecoviruses. In doing so, we aimed to
230 work around a major data limitation: fewer than 20 sarbecovirus hosts were recorded in our dataset, a
231 sample size that would preclude most modeling approaches.

232
233 For both presumed bat host species of sarbecoviruses, the most probable viral sharing hosts were again
234 within the Laurasiatheria. Although bats—especially rhinolophids—unsurprisingly assumed the top
235 predictions given phylogenetic affinity with known hosts (Supplemental Table 3, Supplemental Figure 3),
236 several notable patterns emerged in the rankings of other mammals. Pangolins (Pholidota) were
237 disproportionately likely to share viruses with *R. affinis* and *R. malayanus* (Supplemental Figure 4); the
238 Sunda pangolin (*Manis javanica*) and Chinese pangolin (*M. pentadactyla*) were in the top 20 predictions
239 for both reservoir species (Supplemental Table 4). This result is promising given the much-discussed
240 discovery of SARS-like betacoronaviruses in *M. javanica*¹⁸. The Viverridae were also disproportionately
241 well-represented in the top predictions (Supplemental Figure 5), most notably the masked palm civet
242 (*Paguma larvata*), which was identified as an intermediate host of SARS-CoV^{49,50} (Supplemental Table
243 4). Our virus sharing model thus captured historic patterns of betacoronavirus spillover and predicted
244 cross-species transmission from only the phylogeography of east Asian mammals. The opportunity for
245 inter-species contact depends on both animal behavior and habitat selection, which could be better studied
246 in natural and anthropogenic habitats for these clades⁵¹. However, over evolutionary timescales, our
247 results suggest recombination and diversification of novel betacoronaviruses should be expected among
248 such taxa in nature, without necessarily requiring further contact via the wildlife trade. For example,
249 recent observations in Gabon demonstrate cohabitation between pangolins, bats, and various other
250 mammals (e.g., rodents) in burrows⁵².

251
252 Moreover, these findings lend credibility to other predictions of SARS-CoV-2 sharing patterns and host
253 susceptibility. Many of the model's top predictions were mustelids (i.e., ferrets and weasels), and the most
254 likely viral sharing partner for both *Rhinolophus* species was the hog badger (*Arctonyx collaris*;
255 Supplemental Table 4). Taken together with reports of SARS-CoV-2 spread in mink farms⁵³, these results
256 highlight the relatively unexplored potential for mustelids to serve as betacoronavirus hosts and as models
257 for infectivity studies⁵⁴. Similarly, identification of several deer and Old World monkey taxa as high-
258 probability hosts in our clade-based analysis (Supplemental Figure 3) meshes with the observation of high
259 binding of SARS-CoV-2 to ACE2 receptors in cervid deer and primates⁴⁷. Felids (especially leopards)
260 also ranked relatively high in our viral sharing predictions (Supplemental Table 4, Supplemental Figure
261 5), which is of particular interest given reports of SARS-CoV-2 susceptibility among cats⁵⁵. However, we
262 caution that this model was the only approach in our ensemble that could generate out-of-sample
263 prediction across mammals, and therefore its predictions lacked confirmation (and filtering of potential
264 spurious results) by other models that were designed and implemented independently.

265 266 **Discussion**

267
268 Several limitations apply to our work, most notably the difficulty of empirically verifying predictions.
269 Although some virological studies have incidentally tested specific hypotheses (e.g., filovirus models and

270 bat surveys^{27,56}, henipavirus models and experimental infections^{23,57}), model-based predictions are nearly
271 never subject to systematic verification or post-hoc efforts to identify and correct spurious results. Greater
272 dialogue between modelers and empiricists is necessary to systematically confront the growing set of
273 predicted host-virus associations with experimental validation or field observation.

274
275 Already, we have identified roughly a half-dozen new bat reservoirs of betacoronaviruses that our study
276 correctly predicted. *Scotophilus heathii*, *Hipposideros larvatus*, and *Pteropus lylei*, all highly predicted
277 bat species in our out-of-sample rankings, have been reported positive for betacoronaviruses in the
278 literature^{43,58}; however, resulting sequences were not annotated to genus level in GenBank. More recently,
279 the release of 630 novel bat coronavirus sequences from China included four additional bat hosts
280 (*Hipposideros pomona*, *Scotophilus kuhlii*, *Myotis pequinus*, and *M. horsfieldii*)⁵⁹. All but *M. horsfieldii*
281 were correctly identified by the ensemble, and three models correctly identified all seven verification
282 hosts (Trait-based 1 and 3 and Network-based 1). These results support the idea that our models identified
283 relevant targets correctly but also highlight an evident limitation of the workflow. Whereas an automated
284 approach was the ideal method to systematically compile over 30,000 samples on the timescales
285 commensurate with ongoing efforts to trace SARS-CoV-2 in wildlife, we suggest this discrepancy
286 highlights the need for careful virological work downstream at every stage of the modeling process,
287 including the development of hybrid manual-automated data pipelines.

288
289 Additionally, overcoming underlying model biases that are driven by historical sampling regimes will
290 require coordinated efforts in field study design. Bat sampling for betacoronaviruses has prioritized viral
291 discovery^{39,40,60–62}, but limitations in the spatial and temporal scale (and replication) of field sampling
292 have likely created fundamental gaps in our understanding of infection dynamics in bat populations²⁴.
293 Limited longitudinal sampling of wild bats suggests betacoronavirus detection is sporadic over time and
294 space^{58,63}, implying strong seasonality in virus shedding pulses⁶⁴. Carefully tailored spatial and temporal
295 sampling efforts for priority taxa identified here, within the *Rhinolophus* genus or other high-prediction
296 bat clades, will be key to identifying the environmental drivers of betacoronavirus shedding from wild
297 bats and possible opportunities for contact between bats, intermediate hosts, and humans.

298
299 Future field studies will undoubtedly be important to understand viral dynamics in bats but are inherently
300 costly and labor-intensive. These efforts are particularly challenging during a pandemic, as many
301 scientific operations have been restricted including field studies of bats in some regions to limit possible
302 viral spillback from humans. However, various alternative efforts could both advance basic virology and
303 allow testing model predictions. General open access to viral association records, including GenBank
304 accessions and the upcoming release of the USAID PREDICT program's data, could answer open
305 questions and allow updates to our sampling prioritization (including potentially modeling at subgenus
306 level, with greater data availability). Appropriately preserved museum specimens and historical
307 collections, from diverse research programs, also offer key opportunities to retrospectively screen samples
308 from bats and other mammals for betacoronaviruses and to enhance our understanding of complex host-
309 virus interactions⁶⁵. Large-scale research networks, such as GBatNet (Global Union of Bat Diversity
310 Networks; <https://gbatnet.blogspot.com>) and its member networks, could provide diverse samples and
311 ensure proper partnerships and equitable access and benefit sharing of knowledge across countries^{66,67}.
312 Whole-genome sequencing through initiatives such as the Bat1K Project (<https://bat1k.ucd.ie>) would
313 facilitate fundamental and applied insights into the immunological pathways by which bats can harbor

314 many virulent viruses without displaying disease^{68,69}. Additionally, targeted sequencing efforts could also
315 identify endogenized viral elements in bat genomes, shedding light on the diversity of bat viruses and the
316 evolution of bat immune systems^{70,71}.

317
318 To expedite such work, we have made our binary predictions of host-virus associations for all eight
319 models and all 1,000+ bat species publicly available (Supplementary File 1). Such results are provided
320 both in the spirit of open science and with the hope that future viral detection, isolation, or experimental
321 studies might confirm some of these predictions or rule out others⁵⁷. In ongoing collaborative efforts, we
322 aim to consolidate results from field studies that address these predictions (e.g., serosurveys) and to track
323 Genbank submissions to expand the known list of betacoronavirus hosts. In several years, we intend to
324 revisit these predictions as a post-hoc test of model validation, which would represent the first effort to
325 test the performance of such models and assess their contribution to basic science and to pandemic
326 preparedness.

327
328 It is crucial that our predictions be interpreted as a set of hypotheses about potential host-virus
329 compatibility rather than strong evidence that a particular mammal species is a true reservoir for
330 betacoronaviruses. In particular, susceptibility is only one aspect of host competence^{22,72}, which
331 encompasses the diverse genetic and immunological processes that mediate within-host responses
332 following exposure⁷³. SARS-CoV-2 in particular may have a broad host range⁴⁷, given hypothesized
333 compatibility with the ACE2 receptor in many mammal species, but this only adds to the extreme caution
334 with which any data should be used to implicate a potential wildlife reservoir of the virus, given that rapid
335 interpretation of inconclusive molecular evidence has likely already generated spurious reservoir
336 identifications^{74,75}. Future efforts to isolate live virus from wildlife or to experimentally show viral
337 replication would more robustly test whether predicted host species actually play a role in betacoronavirus
338 maintenance in wildlife⁵⁷.

339
340 Without direct lines of virological evidence, we note that our sampling prioritization scheme also does not
341 implicate any given mammal species in SARS-CoV-2 transmission to humans. Care should be taken to
342 communicate this, especially given the potential consequences of miscommunication for wildlife
343 conservation. The bat research community in particular has expressed concern that negative framing of
344 bats as the source of SARS-CoV-2 will impact public and governmental attitudes toward bat
345 conservation⁷⁶. In zoonotic virus research on bats, studies often over-emphasize human disease risks⁷⁷ and
346 rarely mention ecosystem services provided by these animals⁷⁸. Skewed communication can fuel
347 negative responses against bats, including indiscriminate culling (i.e., reduction of populations by
348 selective slaughter)⁷⁹, which has already occurred in response to COVID-19 even outside of Asia (where
349 spillover occurred)⁸⁰.

350
351 To minimize potential unintended negative impacts for bat conservation, public health and conservation
352 responses should act in accordance with substantial evidence suggesting that culling has numerous
353 negative consequences, not only threatening population viability of threatened bat species in shared
354 roosts⁸¹ but also possibly increasing viral transmission within the very species that are targeted^{82,83}.
355 Instead, bat conservation programs and long-term ecological studies are necessary to help researchers
356 understand viral ecology and find sustainable solutions for humans to live safely with wildlife. From
357 another perspective, policy solutions aimed at limiting human-animal contact could potentially prevent

358 virus establishment in novel species (e.g., as observed in mink farms⁵³), especially in wildlife that may
359 already face conservation challenges (e.g., North American bats threatened by an emerging disease,
360 white-nose syndrome^{79,84}). At least two bat species that can be infected by the fungal pathogen (*Eptesicus*
361 *fuscus* and *Tadarida brasiliensis*) are in our list of the 239 bat species most likely to be betacoronavirus
362 hosts.

363

364 Substantial investments are already being planned to trace the wildlife origins of SARS-CoV-2. However,
365 the intermediate progenitor virus may never be isolated from samples contemporaneous with spillover,
366 and it may no longer be circulating in wildlife. MERS-CoV circulates continuously in camels⁸⁵ and
367 SARS-CoV persisted in civets long enough to seed secondary outbreaks^{49,50}, but the limited description of
368 Pangolin-CoV symptoms suggests high mortality, potentially indicating a more transient epizootic such
369 as Ebola die-offs in red river hogs (*Potamochoerus porcus*)¹⁸. In lieu of concrete data, our study provides
370 no additional evidence implicating any particular species—or any particular pathway of spillover (e.g.,
371 wildlife trade, consumption of hunted animals)—as more or less likely. No specific scenario can be
372 confirmed or rigorously interrogated by ecological models, and we explicitly warn against
373 misinterpretation or misuse of our findings as evidence for adjacent policy decisions. Although policies
374 that focus on particular potential reservoir species or target human-wildlife contact could reduce future
375 spillovers, they will have a negligible bearing on the ongoing pandemic, as SARS-CoV-2 is highly
376 transmissible within humans (e.g., unlike MERS-CoV or other zoonoses that are sustained in people by
377 constant reintroduction). SARS-CoV-2 is likely to remain circulating in human populations at least until a
378 vaccine is developed (if not much longer, given current challenges to vaccine production, allocation,
379 governance, and uptake), regardless of immediate actions regarding wildlife. COVID-19 response must
380 be informed by the best consensus evidence available and prioritize solutions that address immediate
381 reduction of transmission through public health and policy channels. Meanwhile, we hope our proposed
382 wildlife sampling priorities will help increase the odds of preventing the future emergence of novel
383 betacoronaviruses.

384

385 **Acknowledgements**

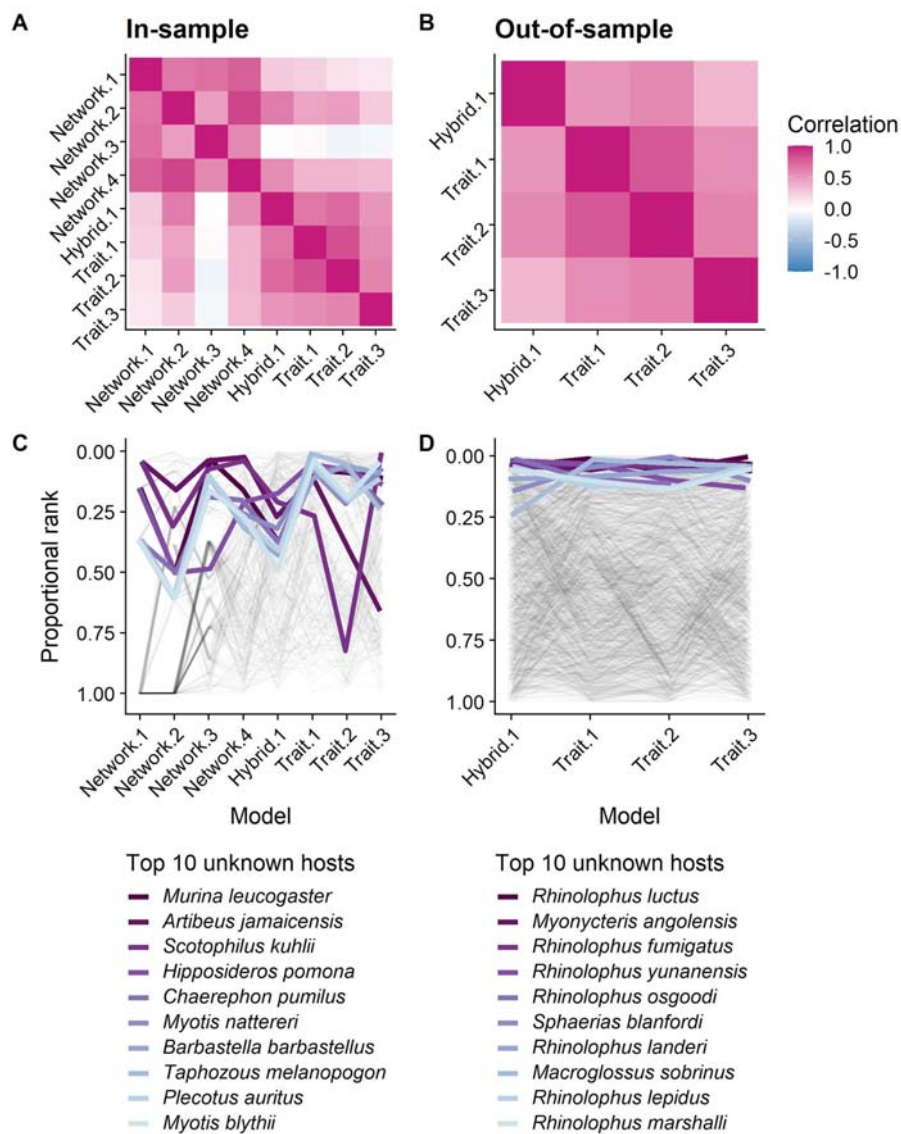
386 We thank Heather Wells for generously sharing thoughtful comments and code. The VERENA
387 consortium is supported by NSF BII 2021909 and by L’Institut de Valorisation de Données (IVADO)
388 through Université de Montreal. DJB was supported by an appointment to the Intelligence Community
389 Postdoctoral Research Fellowship Program at Indiana University, administered by Oak Ridge Institute for
390 Science and Education through an interagency agreement between the U.S. Department of Energy and the
391 Office of the Director of National Intelligence. MS was supported by the Research Foundation - Flanders
392 (FWO17/PDO/067) and the Flemish Government under the “Onderzoeksprogramma Artificiële
393 Intelligentie (AI) Vlaanderen” program.

394

395
396
397
398
399
400
401
402
403
404
405
406

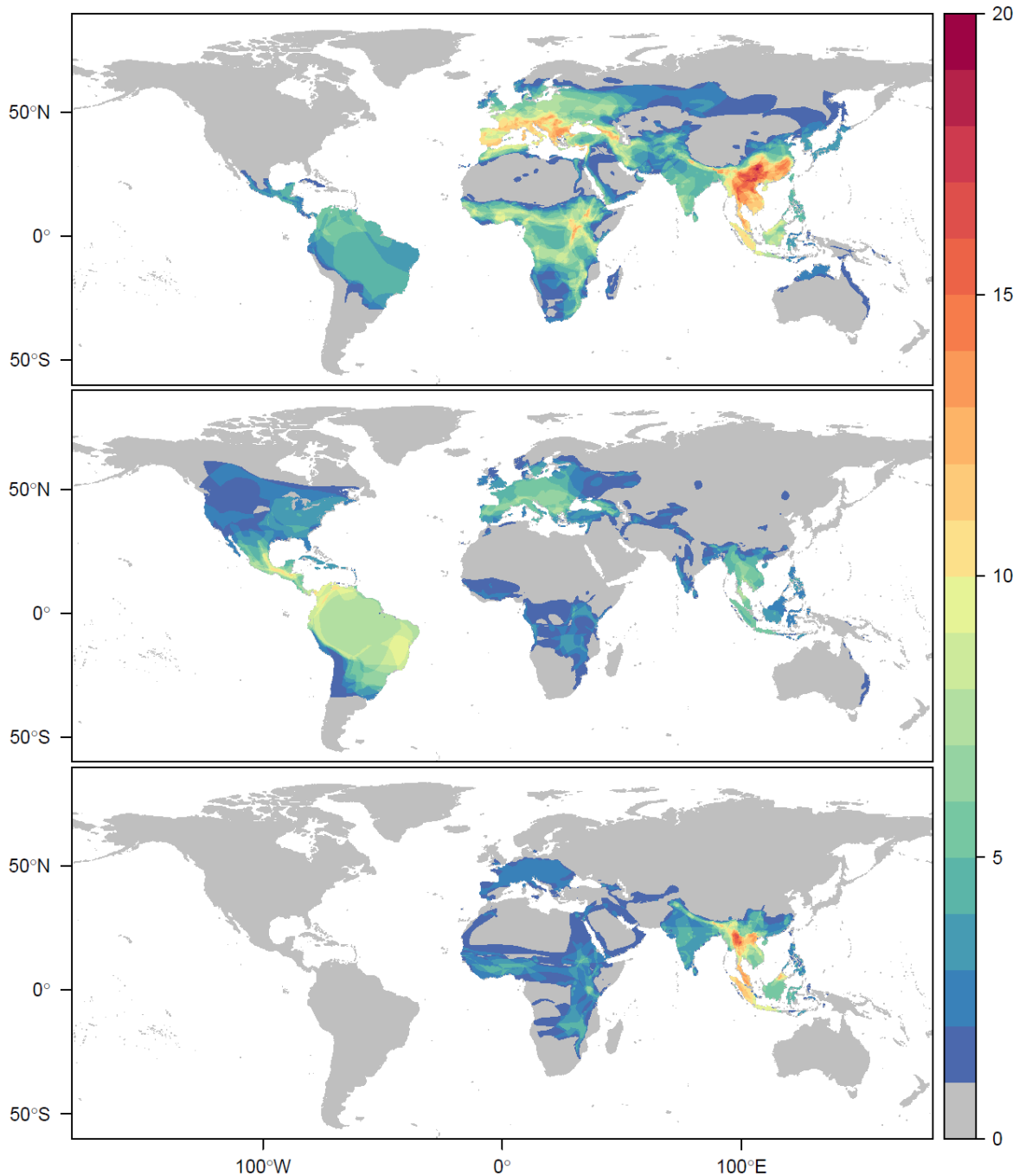
Figures

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.

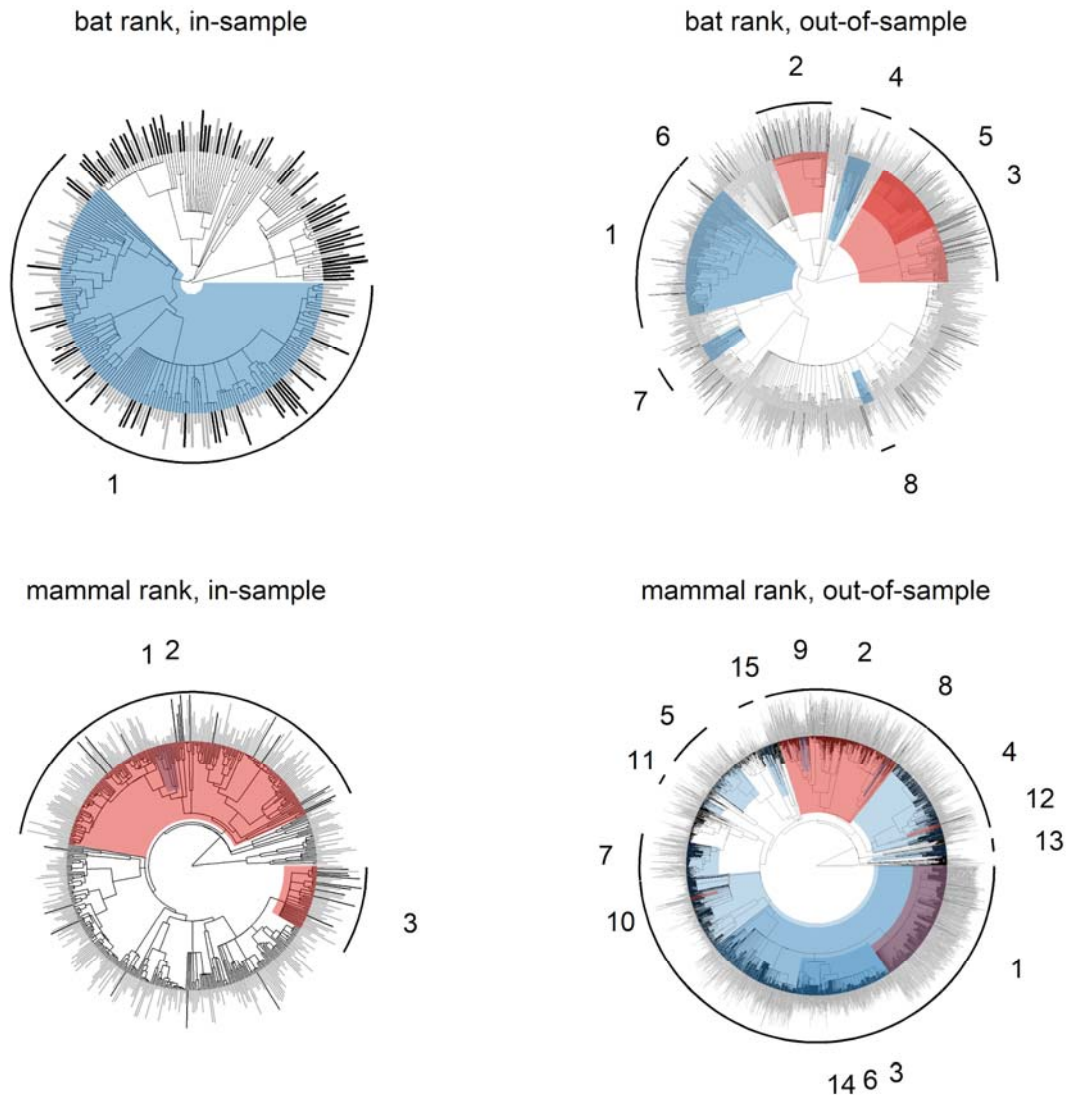


407
408

409 **Figure 2. Species richness of known and suspected betacoronavirus bat hosts.** Known hosts of
410 betacoronaviruses (*top*) are found worldwide, but particularly in southern Asia and southern Europe. The
411 top 50 predicted bat hosts with viral association records (*middle*) are mostly Neotropical, including
412 several species of vampire bats. In contrast, the top 50 *de novo* bat host predictions based on phylogeny
413 and ecological traits (*bottom*) are mostly clustered in Myanmar, Vietnam, and southern China, with none
414 in the Neotropics or North America.
415



417 **Figure 3. Phylogenetic distribution of predicted bat and mammal hosts of betacoronaviruses.** Bar
418 height indicates mean predicted rank across the model ensemble (higher values = lower proportional rank
419 score, more likely to be a host) and black indicates known betacoronavirus hosts. Colored regions indicate
420 clades identified by phylogenetic factorization as significantly different in their predicted rank compared
421 to the paraphyletic remainder; those clades more likely to contain a host are shown in red, whereas those
422 less likely to contain a host are shown in blue. Results are displayed for bats and all mammals separately,
423 stratified by in- and out-of-sample predictions. Numbers reference clade names, species richness, and
424 mean predicted ranks as described in Supplemental Tables 1 and 2.
425



426
427

428
429

430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473

Methods.

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

Host-Virus Association Data

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

Data from GenBank were merged with the Host-Pathogen Phylogeny Project (HP3) dataset³⁰. The HP3 dataset consists of 2,805 associations between 754 mammal hosts and 586 virus species, compiled from the International Committee on Taxonomy of Viruses (ICTV) database, and manually cleaned over a period of five years. Data collection on HP3 began in 2010 and has been static since 2017, but it still represents the most complete dataset on the mammal virome published with a high standard of data documentation. Several recent studies have used the HP3 dataset to produce statistical models of viral sharing or zoonotic potential^{29,48,86}, making it a comparable reference for a multi-model ensemble study.

Because of naming inconsistencies both within GenBank and between the two datasets (HP3 and GenBank), we used a two-step pipeline for taxonomic reconciliation. Viral names were matched to the ICTV 2019 master species list, up to the sub-genus level. Host species names were matched against GBIF using their species API with an automated Julia script, and processed to a fully cleaned set of names. This led to an harmonized dataset representing a global list of mammal-virus associations, from which the bat-coronavirus data can be extracted for downstream and specific modeling efforts. Because the HP3 dataset used an older version of the ICTV master list, and because not all host names in the GenBank metadata could be matched by the GBIF species API (or could be solved unambiguously to the species level), some host-virus interactions were lost; this reinforces the need to careful data curation of taxonomic metadata if they are to enable and support predictive pipelines.

474 **Predictor Data**

475

476 *Phylogeny*

477

478 We used a supertree of extant mammals to unify modeling approaches incorporating host phylogeny³¹.
479 Although more recent mammal supertrees exist, we used this particular phylogeny for consistency with
480 trait datasets and several of the modeling frameworks included in our ensemble. We manually matched
481 select bat species names between our edge list and this particular phylogeny. This included reverting any
482 *Dermanura* to their former *Artibeus* designation (i.e., *D. phaeotis*, *D. cinerea*, *D. tolteca*)⁸⁷, switching
483 *Tadarida* species to either *Mops* or *Chaerephon* species (i.e., *Tadarida condylura* to *Mops condylurus*,
484 *Tadarida plicata* to *Chaerephon plicatus*, *Tadarida pumila* to *Chaerephon pumilus*)⁸⁸, and renaming
485 *Myotis pilosus* to the more recent *Myotis ricketti*. *Chaerephon pusillus* was considered its own species but
486 is now synonymous with *Chaerephon pumilus*⁸⁸. Minor discrepancies between virus data and our
487 phylogeny were also corrected (*Hipposideros commersonii* to *Hipposideros commersoni* [although more
488 recently changed to *Macronycteris commersoni*], *Rhinolophus hildebrandti* to *Rhinolophus hildebrandtii*,
489 *Neoromicia nana* to *Neoromicia nanus*). In other cases, some recently revised genera in our edge list were
490 modified to match former genera in the mammal supertree: *Parastrellus hesperus* to *Pipistrellus*
491 *hesperus*, and *Perimyotis subflavus* to *Pipistrellus subflavus*⁸⁹. Lastly, some names in our edge list
492 missing from the mammal supertree represent former subspecies being raised to full species rank, and
493 names were reverted accordingly: *Artibeus planirostris* to *Artibeus jamaicensis*, *Miniopterus fuliginosus*
494 to *Miniopterus schreibersii*, *Triaenops afer* to *Triaenops persicus*, and *Carollia sowelli* to *Carollia*
495 *brevicauda*. Although we recognize that these are each now recognized as distinct species, in all cases our
496 synonymized names are thought to be either sister taxa or very closely related.

497

498 *Ecological traits*

499

500 We used a previously published dataset of 63 ecological traits describing the morphology, life history,
501 biogeography, and diet of 1,116 bat species. These data are drawn from a combination of PanTHERIA³²,
502 EltonTraits³³, and the IUCN Red List range maps, and were previously cleaned in a study producing
503 predictions of bat reservoirs of filoviruses²⁷. Four redundant variables (two for human population density,
504 mean potential evapotranspiration in range, and body mass) were eliminated prior to analyses, favoring
505 variables with higher completeness.

506

507 *Correction for sampling bias*

508

509 To correct for sampling bias, in the style of several previous studies^{30,86}, we used the number of peer-
510 reviewed citations available on a given host as a measure of scientific sampling effort. We used the R
511 package *easyPubMed* to scrape the number of citations in PubMed returned when searching each of the
512 1,116 bat names in the trait data on April 10, 2020.

513

514 **Modeling Approaches**

515

516 Our team produced an ensemble of eight statistical models (Supplemental Tables 5 and 6), and applied
517 them to generate a predictive set of eight models for bats and five for other mammals. Four use a

518 network-theoretic component (k -nearest neighbors, linear filtering, trait-free plug-and-play, and scaled
519 phylogeny), while three primarily used ecological traits as predictors (boosted regression trees, Bayesian
520 additive regression trees, and neutral phylogeographic). A final hybrid model used a combination
521 approach with a network backbone but incorporating the phylogeny and bat trait data.

522

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

529

530 *Network model 1: k-nearest neighbors recommender*

531

532 We follow the methodology previously developed for the recommendation of species feeding
533 interactions⁹⁰. This method builds a recommender system internally based on the k -nearest neighbor
534 (k NN) algorithm, under which candidate hosts are recommended for a virus from a pool constituted by
535 the hosts of the k viruses with which it has the greatest overlap. Overlap (host sharing) is measured using
536 Tanimoto similarity, which is the cardinality of the intersection of two sets divided by the cardinality of
537 their union. To obtain the pairwise similarity between two viruses, this divides the number of shared hosts
538 by the cumulative number of hosts. The k NN of a virus are the k other viruses with which it has the
539 highest Tanimoto similarity.

540

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

547

548 This model was run two times: first, by measuring the similarity of viruses, and recommending hosts;
549 second, by measuring the similarity of hosts, and recommending viruses. In all cases, only results for
550 betacoronaviruses are reported.

551

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

559

560 *Network-based model 2: Linear filter recommender*

561

562 Following Stock *et al.*⁹¹, we used a previously developed linear filter to infer potential missing
563 interactions. This recommender system assumes that networks tend to be self-similar, and use this
564 information to generate a score for an un-observed interaction that is a linear combination of the status of
565 the interaction (relative weight of 1/4), relative degree of host and virus, and of the observed connectance
566 of the network (all with relative weights of 1); as we are concerned with ranking interactions as opposed
567 to examining the absolute value of the score, the penalization coefficient associated to the interaction
568 being presumed absent could be omitted with no change in the ranking, but has been set to a low value
569 instead. The scores returned by the linear filter are not directly related to the probability of the interaction
570 existing in this context, but higher scores still indicate interactions that are more likely to exist. Indeed,
571 known hosts of betacoronavirus typically scored higher.

572
573 We used the zero-one-out approach to assess the performance of this model on the entire datasets. In all
574 cases, non-interactions ranked lower than positive interactions even when entirely removing the
575 penalization coefficient from the linear filter parameters, which suggests that the network structure
576 (degree and connectance) is capturing a lot of information as to which species can interact.

577
578 *Network-based model 3: Plug and play*

579
580 For network problems, the “plug and play” model is a statistical approach that formulates Bayes’ theorem
581 for link prediction around the conditional density of traits of known associations compared to traits of
582 every possible association in a network. The conditional density function is measured by using non-
583 parametric kernel density estimators (implemented with the R package *np*), and the conditional ratio
584 between them is used to estimate link “suitability”, a scale-free ratio. Compared to other machine learning
585 methods that fit training data iteratively, plug and play is comparatively simple, and directly infers the
586 most likely extensions of observed patterns in data. The plug and play was originally developed to
587 forecast missing links in host-parasite networks³⁶, but has since been used to model species distributions⁹²
588 and predict the global spread of human infectious diseases⁹³. We used this model here to estimate
589 suitability of host-virus interactions by first modeling the entire estimated network of host-virus
590 interaction suitability, and ranking hosts that are not infected by betacoronaviruses by their estimated
591 suitability for betacoronaviruses.

592
593 The “plug and play” model is trained using either matched pairs of host and pathogen ecological,
594 morphological, or phylogenetic traits³⁶, or by using a latent approach⁹³ which considers the mean
595 similarity of pathogens in their host ranges and the mean similarity of hosts in their pathogen
596 communities as ‘traits’. We decided to use the latent approach, as host trait data was far more available
597 than viral trait data. Further, the taxonomic scale considered for host (species) and virus (genus) differed,
598 making the resolution of potential trait data different enough to potentially confound trait-based
599 approaches in this modeling framework.

600
601 Relative suitability of a host-virus association, as estimated by the “plug and play” model, is formulated
602 as a density ratio estimation problem. The suitability of a host-virus association is quantified as the
603 quotient of the distribution of latent trait values when an association was recorded over the distribution of
604 all the latent trait values. As an attempt to control for sampling effort of mammal and bat host species, we
605 included PubMed citation counts for host species (as described above) in the estimation of host-virus

606 suitability. We explored host-pathogen suitability using the entire mammal-virus associations dataset, to
607 maximize the available information on the network's structure, and ranked host-pathogen pairs by their
608 relative suitability value. From the final predictions, we subset out bat-specific predictions. When
609 predicting, we set citation counts to the mean of training data, as a sampling bias correction.

610

611 *Network-based model 4: Scaled-phylogeny*

612

613 We apply the network-based conditional model of Elmasri *et al.*⁹⁴ for predicting missing links in bipartite
614 ecological networks. The full model combines a hierarchical Bayesian latent score framework which
615 accounts for the number of interactions per taxon, and a dependency among hosts based on evolutionary
616 distances. To predict links based on evolutionary distance, the probability of a host-parasite interaction is
617 taken as the sum of evolutionary distances to the documented hosts of that parasite. This allocates higher
618 probabilities when a few closely related hosts, or many distantly related hosts interact with a parasite. In
619 this way phylogenetic distances are combined with individual affinity parameters per taxa to model the
620 conditional probability of an interaction.

621

622 In ecological studies, it is common to use time-scaled phylogenies to quantify evolutionary distance
623 among species⁹⁵. We may use these fixed evolutionary distances for link prediction, but parasite taxa are
624 known to be more or less constrained by phylogenetic distances among hosts⁹⁶. Further, phylogenies are
625 hypotheses about evolutionary relationships and have uncertainties in the topology and relative distances
626 among species⁹⁷. Rather than treating phylogenetic distances as fixed, Elmasri *et al.*⁹⁴ re-scale the
627 phylogeny by applying a macroevolutionary model of trait evolution. While any evolutionary model that
628 re-scales the covariance matrix may be used, we use the early-burst model, which allows evolutionary
629 change to accelerate or decelerate through time⁹⁸. This different emphasis to be placed on deep versus
630 recent host divergences when predicting links.

631

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

640

641 To assess predictive performance, we attempted to predict the held out interactions, and calculated AUC
642 scores by thresholding predicted probabilities per fold, and taking an average across the 10 folds. In
643 addition to AUC, we also assessed the model based on the percent of documented interactions accurately
644 recovered. For the bat-viral genera data the full model resulted in an average AUC of 0.82 and recovered
645 an average of 90.1% of held out interactions, while the phylogeny-only model showed increased AUC
646 (0.86), but a decreased proportion of held-out interactions recovered (84.5%). Interestingly, the models
647 for bat-virus genera associations had marginally worse predictive performance compared to the same
648 models run on the larger network of mammal-virus associations (full model: AUC 0.88, 84.4% positive
649 interactions recovered; phylogeny-only model: AUC: 0.88, 88.8% positive interactions recovered),

650 indicating that predicting bat-betacoronavirus associations may benefit from including data on non-bat
651 hosts. The models also estimated the scaling parameter (η) of the early-burst model to be positive
652 (average $\eta=7.92$ for the full model run on the bat subset), indicating accelerating evolution compared to
653 the input tree (Supplemental Figure 6). This means that recent divergences are given more weight than
654 deeper ones for determining bat-viral genera associations, which is consistent with recent work on viral
655 sharing^{48,99}.

656

657 *Trait-based model 1: Boosted regression trees*

658

659 Previous work has been highly successful in predicting zoonotic reservoirs using a combination of
660 taxonomic, ecological, and geographic traits as predictors. This approach has been previously used to
661 identify wildlife hosts of filoviruses^{27,100}, flaviviruses^{28,101}, henipaviruses²³, *Borrelia burgdorferi*²⁶, to
662 predict mosquito vectors of flaviviruses¹⁰², and to predict rodent reservoirs and tick vectors of zoonotic
663 viruses^{37,103}. These approaches treat the presence of a specific virus (or genus of viruses) or a zoonotic
664 pathogen as an outcome variable, with negative values given for species not known to be hosts
665 (pseudoabsences), and use machine learning to identify the characteristics that predispose animals to
666 hosting pathogens of concern. By predicting the probability that a given pseudoabsence is a false
667 negative, the method can infer potential undetected or undiscovered host species.

668

669 This approach has almost exclusively been implemented using boosted regression trees (BRT), a
670 classification and regression tree (CART) machine learning method that became popular a decade ago for
671 species distribution modeling.¹⁰⁴ Boosted regression trees develop an ensemble of classification trees
672 which iteratively explain the residuals of previous trees, up to a fixed tree depth (usually between 3 and 5
673 splits). The incorporation of boosting allows the model, as it is fit, to progressively better explain poorly-
674 fit cases within training data.

675

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

683

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

688

689 We applied bootstrapping resampling methods to estimate uncertainty, using our tuned model to fit 1000
690 replicate models. For each model, training sets were assembled by randomly selecting with replacement
691 79 bat-coronavirus associations from the set of reported bat hosts and 79 pseudoabsences. Trained models
692 were used to generate relative influence coefficients for trait predictors and coronavirus host probabilities
693 across all bat species. Partial dependence plots display relative influence coefficients and bootstrapped

694 confidence intervals for the top ten contributing trait predictors. The medians of host probabilities were
695 ranked and used to identify the top ten candidate host species. When predicting, we set citation counts to
696 the mean of training data, as a sampling bias correction.

697

698 *Trait-based model 2: Bayesian additive regression trees*

699

700 A similar workflow to trait-based model 1 was implemented using Bayesian additive regression trees
701 (BART), an emerging machine learning tool that has similarities to more popular methods like random
702 forests and boosted regression trees. BART adds several layers of methodological innovation, and
703 performs well in bakeoffs with other advanced machine learning methods. Several features make BART
704 very convenient for modeling projects like these, including several easy-to-use implementations in R
705 packages, built-in capacity to impute and predict on missing data, and easy construction of variable
706 importance and partial dependence plots.

707

708 Like other classification and regression tree methods, BART assigns the probability of a binary outcome
709 variable by developing a set of classification trees - in this case, a sum-of-trees model - that split data
710 (“branches”) and assign values to terminal nodes (“leaves”). Whereas other similar methods generate
711 uncertainty by adjusting data (e.g. random forests bootstrap training data and fit a tree to each bootstrap;
712 boosted regression trees are usually implemented with iterated training-test splits to generate confidence
713 intervals), BART generates uncertainty using an MCMC process. An initial sum-of-trees model is fit to
714 the entire dataset, and then rulesets are adjusted in a limited and stochastic set of ways (e.g., adding a
715 split; switching two internal nodes), with the sum-of-trees model backfit to each change. After a burn-in
716 period, the cumulative set of sum-of-trees models is treated as a posterior distribution. This has some
717 advantages over other methods, like boosted regression trees or random forests. In particular, posterior
718 width directly measures model uncertainty (rather than approximating it by permuting training data), and
719 a single model can be run (instead of an ensemble trained on smaller subsets of training data), allowing
720 the model to use the full training dataset all at once.¹⁰⁵

721

722 Unlike many Bayesian machine learning methods, BART is easily implemented out-of-the-box, due to a
723 limited set of customization needs. Three main priors control the fitting process: one usually-uniform
724 prior on variable importance, one two-parameter negative power distribution on tree depth (preventing
725 overfitting), and an inverse chi-squared distribution on residual variance. A set of well-performing priors
726 from the original BART study¹⁰⁶ are widely used across R implementations for out-of-the-box settings,
727 but can be further adjusted relative to modeling needs. In this study, we implemented BART models
728 using a Dirichlet prior for variable importance (DART), a specification that is designed for situations with
729 high dimensionality data that probably reflects a small number of true informative predictors. This often
730 produces a much more reduced model without going through a stepwise variable selection process, which
731 can be slow and very subject to stochasticity.¹⁰⁵

732

733 We implemented this approach using the *BART* package in R, using the bat-virus association dataset to
734 generate an outcome variable, and the bat traits dataset as predictors. BART models were implemented
735 with 200 trees and 10,000 posterior draws, using every trait feature that was at least 50% complete and <
736 97% homogenous (taken from TBM1).

737

738 We tried four total implementations, based on two decisions: BART uncorrected and corrected for
739 citation counts (BART-u, BART-c), and DART uncorrected and corrected for citation counts (DART-u,
740 DART-c). All four models performed well, with little variation in predictive power measured by the area
741 under the receiver operator curve calculated on training data (BART-u: AUC = 0.93; BART-c: AUC =
742 0.93; DART-u: AUC = 0.93; DART-c: 0.90; Supplemental Figure 7). Across all models, spatial variables
743 had a high importance, including some regionalization (extent of range) and some variables capturing
744 larger geographic range sizes, as did a diet of invertebrates (pulling out the phylogenetic signal of
745 insectivorous bats; Supplemental Figure 8).

746
747 All models identified a number of “false negative” hosts that would be suitable based on a 10% false
748 negative classification threshold for known betacoronavirus hosts (implemented with the R package
749 ‘PresenceAbsence’). BART-u identified 217 missing hosts, BART-c identified 279 missing hosts, DART-
750 u identified 222 missing hosts, and DART-c identified 384 missing hosts, suggesting that this model most
751 penalized overfitting as intended. As a result, we considered this model the most rigorous and powerful
752 for inference, and used DART-c in the final model ensemble. We predicted across all 1,040 bats without
753 recorded betacoronavirus associations, and ranked predicted probability. When predicting, we set citation
754 counts to the mean of training data, as a sampling bias correction.

755

756 *Trait-based model 3: Phylogeographic neutral model*

757
758 We used a previously published pairwise viral sharing model⁴⁸ to predict potential betacoronavirus hosts
759 based on the sharing patterns of known hosts in a published dataset³⁰. We used a generalised additive
760 mixed model (GAMM), which was fitted in the first half of 2019 using the *mgcv* package, with pairwise
761 binary viral sharing (0/1 denoting if a species shares at least one virus) as a response variable.
762 Explanatory variables include pairwise proportional phylogenetic distance and geographic range overlap
763 (taken from the IUCN species ranges), with a multi-membership random effect to control for species-
764 level sampling biases. The model was then used to predict the probability that a given species pair share
765 at least one virus across 4196 placental mammals with available data, producing a predicted viral sharing
766 network that recapitulates a number of known macroecological patterns, as well as predicting reservoir
767 hosts with surprising accuracy⁴⁸. Subsetting this predicted sharing matrix, we listed the rank order of hosts
768 most likely to share with all known betacoronavirus hosts in our datasets.

769

770 *Rhinolophus-specific implementation of trait-based model 3*

771
772 We then repeated this process with sharing patterns of *Rhinolophus affinis* and *R. malayanus* specifically.
773 Given the strong phylogenetic effect, the top 139 predictions were bat species: predominantly
774 rhinolophids and hipposiderids. The top 20 predictions for both *R. malayanus* and *R. affinis* are displayed
775 in Supplemental Tables 3 and 4. Notable predictions included the hog badger *Arctonyx collaris*
776 (Carnivora: Mustelidae), which was examined for SARS-CoV antibodies in 2003 and is reported in
777 wildlife markets^{7,107}; a selection of civet cats (Carnivora: Viverridae) including *Viverra* species; the
778 binturong (*Arctitis binturong*); and the masked palm civet (*Paguma larvata*), the latter of which were
779 implicated in the chain of emergence for SARS-CoV^{49,50}; and pangolins (Pholidota: Manidae) including
780 *Manis javanica* and *Manis pentadactyla*, which have been hypothesised to be part of the emergence chain
781 for SARS-CoV-2^{18,19}.

782

783 Alongside these high-ranked species-level predictions, we visually examined how predictions varied
784 across all mammal orders and families using the whole dataset (Supplemental Figure 5). Pangolins
785 (Pholidota), treeshrews (Scandentia), carnivores (Carnivora), hedgehogs (Erinaceomorpha), and even-
786 toed ungulates (Artiodactyla) had high mean predicted probabilities. Investigating family-level sharing
787 probabilities revealed that civets (Viverridae) and mustelids (Mustelidae) were responsible for the high
788 Carnivora probabilities, and mouse deer (Tragulidae) and bovids (Bovidae) were mainly responsible for
789 high probabilities in the Artiodactyla (Supplemental Figure 5).

790

791 *Hybrid model 1: Two-Step Kernel Ridge Regression*

792

793 To make predictions using both the traits and the phylogeny of the hosts, we used Two-Step Kernel Ridge
794 Regression (TSKRR), a kernel-based pairwise learning method^{108,109}. This method can predict interactions
795 between two species, i.e., hosts and viruses, both in-sample and out-of-sample. Conceptually, TSKRR
796 performs nonlinear regression twice: once to generalize to new hosts and once to generalize to new
797 viruses (Supplemental Figure 9A). As a kernel method, it can exploit arbitrary similarity measures to
798 describe the species. At the same time, its conceptual simplicity permits using efficient shortcuts for
799 tuning and specialized cross-validation.

800

801 For the hosts, we again restricted ourselves to bats. First, we inputted missing trait values using the
802 *MissForest* package¹¹⁰. We then selected the relevant traits, standardized them, and computed the squared
803 Euclidean distance. We plugged this distance into a standard radial basis kernel¹¹¹, where we used the
804 median heuristic to set the bandwidth. Similarly, for the bat phylogeny, we computed the distance matrix.
805 We combined this distance with the radial basis kernel, again setting the bandwidth by the median of the
806 distances. Bats were described using a kernel matrix derived from traits, phylogeny, or the average of
807 both (Supplemental Figure 9B). We corrected for sampling bias by including the number of citations per
808 species as a covariate and predicted using the average number of citations for all hosts. For viruses, we
809 used a simple kernel with a “1” if two virus genera belong to the same family, and a “0” otherwise. For all
810 the kernel matrices, we added 0.1 to each element and an additional 0.1 to the diagonal elements. This
811 modification serves as an intercept and it gives the TSKRR the same capacity as the linear filtering
812 method described earlier¹⁰⁸.

813

814 We fitted the models using the R package *xnet*. We assessed model performance in two ways. First, we
815 used leave-one-out for every element of the incidence matrix and computed the AUC over all predicted
816 interactions. Secondly, we used leave-one-out cross-validation on bats(i.e., leaving out a single species,
817 making predictions for all the viruses of that bat, and repeating this for every species). In this setting, we
818 computed the AUC for every virus over all bats and averaged these values (Supplemental Figure 9C). For
819 both evaluation strategies, we removed viruses with no hosts in the dataset, although we kept these for
820 training the model (Supplementary Table 7). Compared to the average performance for the viruses,
821 betacoronaviruses were slightly more difficult to predict. Phylogeny was slightly more informative than
822 traits, although the effect between them was synergistic. As would be expected, it was easier to predict
823 missing interactions than to predict new hosts.

824

825 To assess which variables are most important, we computed feature importance based on random
826 permutations. We randomly reshuffled one or both of the bats' kernel matrices, destroying its information
827 content. Subsequently, we recorded the decrease in AUC for the setting 'interactions' in the model as an
828 indication of relative importance. We also randomly reshuffled each trait separately and monitored the
829 effect on performance. We repeated each random permutation 100 times. Reshuffling the complete kernel
830 matrix with both traits and phylogeny resulted in the largest average drop in AUC of 0.1904
831 (Supplementary Figure 10). Reshuffling only the phylogeny kernel matrix had a much more profound
832 impact than the trait matrix (0.07769 vs. 0.05301). Contributions of traits were smaller by comparison.
833 Citation count was the largest effect (0.02119), with diet breadth (0.00520), population group size
834 (0.003802), and annual birth pulse (0.0022944) as the remaining largest contributions. Removing
835 citations resulted in a stronger importance for these variables but did not change the relative importance.

836

837 **Consensus Methods and Recommendations**

838

839 *Combining and ranking predictions*

840

841 For all eight models predicting bat hosts of betacoronaviruses, and five models predicting mammal hosts
842 of betacoronaviruses, we combined predictions—generated using the same standardized data—into one
843 standardized dataset. All mammal models were trained on data including bats, but predictions were subset
844 to exclude bats to focus on likely intermediate hosts.

845

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

853

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

858

859 For bats, the final ensemble of models spanned a large range of performance on the training data,
860 measured by the area under the receiver operator curve (AUC; Network 1: 0.624; Network 2: 0.987;
861 Network 3: 0.514; Network 4: 0.726; Trait 1: 0.850; Trait 2: 0.902; Trait 3: 0.762; Hybrid 1: 0.924),
862 indicating that it was possible to suitably detect differences in model performance on the full data. The
863 total ensemble of proportional ranks performed medium well (AUC = 0.832). We used known
864 betacoronavirus associations to threshold each model and the ensemble predictions based on a 10%
865 omission threshold (90% sensitivity), and we again found a wide range in the number of predicted
866 undiscovered bat hosts of betacoronaviruses (Network 1: 162 species; Network 2: 1; Network 3: 111;
867 Network 4: 44; Trait 1: 425; Trait 2: 384; Trait 3: 720; Hybrid 1: 181; total ensemble: 239 species). Given

868 concerns about mammal model performance and biological accuracy (see Main Text), we elected not to
869 apply this exercise to mammal hosts at large.

870
871 To visualize the spatial distribution of predicted bat hosts, we used the IUCN Red List database of species
872 geographic distributions. We took the top 50 ranked in-sample predictions and top 50 ranked out-of-
873 sample predictions and combined these range maps to visualize species richness of top predicted hosts
874 (Figure 3).

875
876 *Phylogenetic factorization of ensemble models*

877
878 We used phylogenetic factorization to flexibly identify taxonomic patterns in the consensus proportional
879 rankings of likely hosts of SARS-CoV-2. Phylogenetic factorization is a graph-partitioning algorithm that
880 iteratively partitions a phylogeny in a series of generalized linear models to identify clades at any
881 taxonomic level (e.g., rather than *a priori* comparing strictly among genera or family) that differ in a trait
882 of interest⁴⁵. Using the mammal supertree, we used the *phylofactor* package to partition proportional rank
883 as a Gaussian-distributed variable. We determined the number of significant phylogenetic factors using a
884 Holm's sequentially rejective 5% cutoff for the family-wise error rate. We applied this algorithm across
885 our four final ensemble prediction datasets: in-sample bat ranks, out-of-sample bat ranks, in-sample
886 mammal ranks, and out-of-sample mammal ranks.

887
888 Using network and trait-based models within-sample, we identified only one bat clade with substantially
889 different consensus proportional rankings, the Yangochiroptera ($x_{\square}=0.55$ compared to 0.39 for the
890 remaining bat phylogeny, the Yinpterochiroptera). Out of sample, using only trait-based models, we
891 instead identified eight bat clades with different propensities to include unlikely or likely bat hosts of
892 betacoronaviruses. Subclades of the New World superfamily Noctilionoidea broadly had higher
893 proportional ranks ($x_{\square}=0.73$), indicating lower predicted probability of being hosts, as did the
894 Emballanuridae ($x_{\square}=0.75$). In contrast, both the Rhinolophidae and the Pteropodidae broadly had lower
895 mean ranks ($x_{\square}=0.29$ and $x_{\square}=0.38$).

896
897 Using network models within-sample across non-volant mammals, we identified four clades with
898 different proportional ranks. The largest clade was the Laurasiatheria (Artiodactyla, Perissodactyla,
899 Carnivora, Pholidota, Soricomorpha, and Erinaceomorpha), which had lower proportional ranks (higher
900 risk; $x_{\square}=0.55$). Nested within this clade, the Cetacea had greater proportional ranks ($x_{\square}=0.89$),
901 indicating lower risk. A large subclade of the Murinae (Old World rats and mice) also had lower ranks
902 ($x_{\square}=0.52$). Out of sample, using only the biogeographic viral sharing model, we instead identified 15
903 clades with different proportional ranks. The first clade identified large swaths of the Muridae as having
904 higher risk ($x_{\square}=0.38$) as well as the Laurasiatheria ($x_{\square}=0.50$). Old World primates had weakly lower risk
905 ($x_{\square}=0.65$), as did the Scuridae ($x_{\square}=0.67$). The Cetacea and Pinnipedia both had greater proportional
906 ranks ($x_{\square}=0.89$ and $x_{\square}=0.71$). Old World porcupines (Hystricidae) and the Erinaceidae (Paraechinus,
907 Hemiechinus, Mesechinus, Erinaceus, Atelerix) both had greater risk ($x_{\square}=0.48$ and $x_{\square}=0.39$), while the
908 Afrosoricida had higher ranks ($x_{\square}=0.97$).

909
910 To assess potential discrepancy between taxonomic patterns in model ensemble predictions and those of
911 simply host betacoronavirus status itself, we ran a secondary phylogenetic factorization treating host

912 status as a Bernoulli-distributed variable, with the same procedure applied to determine the number of
913 significant phylogenetic factors. To assess sensitivity of taxonomic patterns to sampling effort, we ran
914 phylogenetic factorization with and without square-root transformed PubMed citations per species as a
915 weighting variable (Supplemental Figure 11).

916
917 Without accounting for study effort, phylogenetic factorization of betacoronavirus host status identified
918 one significant clade across the bat phylogeny, the Yangochiroptera, as having fewer positive species
919 (4.71%) than the paraphyletic remainder (12.12%). When accounting for study effort, however, the single
920 clade identified by phylogenetic factorization changed, with a subclade of the family Pteropodidae (the
921 Rousettinae) having a greater proportion of positive species (28.6%). For non-volant mammals,
922 phylogenetic factorization identified only one clade, the family Camelidae, as having more positive
923 species (75%) than the tree remainder (0.68%).

924
925 *Phylogenetic factorization of Rhinolophidae virus sharing*

926
927 Because phylogenetic patterns in predictions from our viral sharing model could vary across other
928 taxonomic scales beyond order and family, we also used phylogenetic factorization to more flexibly
929 identify host clades with different propensities to share viruses with *R. affinis* and *R. malayanus*. We
930 partitioned rank as a Gaussian-distributed variable and again determined the number of significant
931 phylogenetic factors using Holm's sequentially rejective 5% cutoff.

932
933 Within the Chiroptera, we identified 10 clades with different propensities to share viruses with *R. affinis*
934 and 5 clades with different propensities to share viruses with *R. malayanus*. For both bats, the top clade
935 was the family Rhinolophidae, reinforcing phylogenetic components of the biogeographic model and
936 highlighting the greater likelihood of viral sharing (mean rank $x_{\square}=40$ for *R. affinis*, $x_{\square}=42$ for *R.*
937 *malayanus*). For *R. affinis*, several individual bat species had lower risks of viral sharing (e.g., *Myotis*
938 *leibii*, $x_{\square}=4100$; *Pteropus insularis*, $x_{\square}=3157$; *Nyctimene aello*, $x_{\square}=2497$; *Chaerephon chapini*,
939 $x_{\square}=2497$). The Megadermatidae, Nycteridae, and Hipposideridae (under which the PanTHERIA dataset
940 includes the genus *Rhinonictoris*, although this is now considered a separate family, the
941 Rhinonycteridae¹¹²) collectively had greater likelihood of viral sharing ($x_{\square}=557$), as did the
942 Vespertilionidae ($x_{\square}=704$).

943
944 Across the non-volant mammals, we identified 7 clades with different propensities to share viruses with
945 *R. affinis* and only 1 clade with different propensities to share viruses with *R. malayanus*. For both bat
946 species, the first and primary clade was the Ferungulata (Artiodactyla, Perissodactyla, Carnivora,
947 Pholidota, Soricomorpha, and Erinaceomorpha), which had lower ranks (higher viral sharing; $x_{\square}=2084$).
948 For viral sharing with *R. affinis*, the Sciuridae was more likely to share viruses ($x_{\square}=1948$), as was the
949 Scandentia ($x_{\square}=1416$) and many members of the Colobinae ($x_{\square}=1958$). However, members of the tribe
950 Muntiacini (genera *Elaphodus* and *Muntiacus*) had especially high likelihoods of viral sharing and low
951 rank ($x_{\square}=361$).

952
953 **Data and Code Availability**

954

955 The standardized data on betacoronavirus associations, and all associated predictor data, is available from
956 the VERENA consortium's Github (github.com/viralemergence/virionette). All modeling teams
957 contributed an individual repository with their methods, which are available in the organizational
958 directory (github.com/viralemergence). All code for analysis, and a working reproduction of each
959 authors' contributions, is available from the study repository (github.com/viralemergence/Fresnel).

960
961
962
963
964
965
966
967
968

Supplemental Material

Supplemental 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	0.550	0.393
out	1	Mystacinidae, Noctilionidae, Mormoopidae, Phyllostomidae	161	0.730	0.483
out	2	Rhinolophidae	73	0.295	0.538
out	3	Pteropodidae	173	0.379	0.548
out	4	Mosia, Emballonura, Coleura, Rhynchonycteris, Cyttarops, Diclidurus, Centronycteris, Cormura, Saccopteryx, Balantiopteryx, Peropteryx	31	0.755	0.512
out	5	Eonycteris, Macroglossus, Syconycteris, Notopteryx, Melonycteris, Harpyionycteris, Aproteles, Dobsonia, Styloctenium, Neopteryx, Pteralopex, Acerodon, Pteropus	99	0.458	0.527
out	6	Thyropteridae, Furipteridae, Natalidae	12	0.818	0.517
out	7	Otomops, Promops, Molossus, Eumops	26	0.669	0.517
out	8	Lasiurus	15	0.717	0.518

969
970

971 **Supplemental Table 2. Results of phylogenetic factorization applied to predicted rank probabilities**
 972 **for all mammals.** The number of retained phylogenetic factors (following a 5% family-wise error rate
 973 applied to GLMs), taxa corresponding to those clades, number of species per clade, and mean predicted
 974 rank probabilities for the clade compared to the paraphyletic remainder are shown stratified by models
 975 applied in- and out-of-sample.
 976

Sample	Factor	Taxa	Tips	Clade	Other
in	1	Phocoenidae, Delphinidae, Tursiops, Monodontidae, Physteridae, 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, Lemniscomys, Mastomys, Mus, Pelomys, Niviventer, Otomys, Praomys, Rattus, Vandeleuria	38	0.520	0.627
out	1	Abditomys, Bullimus, Limnomys, Tarsomys, Trypomys, Acomys, Lophuromys, Uranomys, Aethomys, Micaelamys, Anisomys, Chirurmys, 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, Grammomy, Thallomys, Hadromys, Heimyscus, Hylomyscus, Komodomys, Papagomys, Oenomys, Thamnomys, Lemniscomys, 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, Platanthomyidae	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

977

978 **Supplemental Table 2, continued.** (Page 2 of 2)

979

Sample	Factor	Taxa	Tips	Clade	Other
out	6	Abrawayaomys, 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, Lundomys, 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, Physteridae, 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, Dasyproctidae, 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

980

981 **Supplemental Table 3. Predicted high-similarity bat hosts sharing with *Rhinolophus affinis* and *R.***
 982 ***malayanus*.** Species on these lists may be particularly likely to be the ultimate evolutionary origin of
 983 SARS-CoV-2, or a closely-related virus prior to recombination in an intermediate host. Predictions are
 984 made based just on the average viral sharing probability inferred for the two hosts from the
 985 phylogeography model (Trait-based 3). (* Note that the two species have high sharing probabilities with
 986 each other, potentially indicating that efforts to trace the origins of SARS-CoV-2 are already very close to
 987 their target.)
 988

Rhinolophus affinis	Rhinolophus malayanus
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)

989

990 **Supplemental Table 4. Predicted high-similarity non-bat hosts sharing with *Rhinolophus affinis* and**
 991 ***R. malayanus*.** Species on these lists may be particularly suitable as stepping stones for betacoronavirus
 992 transmission from bats into humans, including potentially for SARS-CoV-2 and other SARS-like viruses.
 993 Predictions are made based just on the average viral sharing probability inferred for the two hosts from
 994 the phylogeography model (Trait-based 3). Species' binomial names are included alongside their families.
 995

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

996

997 **Supplemental Table 5. Taxonomic scale of model training data and predictive implementation.**

998 Notes: (1) These models generated predictions of sharing with *Rhinolophus affinis* over all non-human
 999 mammals in the HP3 dataset, then subsetted to bats. (2) In these models, bat-betacoronavirus predictions
 1000 are based on a subset of binary outcomes for known association with betacoronaviruses, without any other
 1001 viruses included.
 1002

Model approach	Training data scale	Bat <i>Betacoronavirus</i> predictions	Mammal-wide <i>Betacoronavirus</i> predictions
Network-based 1 k-Nearest neighbors	Bat-virus	<input type="checkbox"/>	
Network-based 1 k-Nearest neighbors	Mammal-virus		<input type="checkbox"/>
Network-based 2 Linear filter	Bat-virus	<input type="checkbox"/>	
Network-based 2 Linear filter	Mammal-virus		<input type="checkbox"/>
Network-based 3 Plug-and-play	Mammal-virus ¹	<input type="checkbox"/>	<input type="checkbox"/>
Network-based 4 Scaled-phylogeny	Bat-virus	<input type="checkbox"/>	
Network-based 4 Scaled-phylogeny	Mammal-virus		<input type="checkbox"/>
Trait-based 1 Boosted regression trees	Bat-betacoronavirus ²	<input type="checkbox"/>	
Trait-based 2 Bayesian additive regression trees	Bat-betacoronavirus ²	<input type="checkbox"/>	
Trait-based 3 Neutral phylogeographic	Mammal-virus ¹	<input type="checkbox"/>	<input type="checkbox"/>
Hybrid 1 Two-step kernel ridge regression	Bat-betacoronavirus	<input type="checkbox"/>	

1003

1004

1005 **Supplemental Table 6. Data scale of prediction, by method.** Some methods use pseudoabsences to
 1006 expand the scale of prediction but still only analyze existing data, with no out-of-sample inference, while
 1007 others can predict freshly onto new data. (* Training data from the HP3 database uses pseudoabsences,
 1008 but no new ones are generated in this study that modify the model or the bat-virus association dataset)
 1009

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

1010
1011

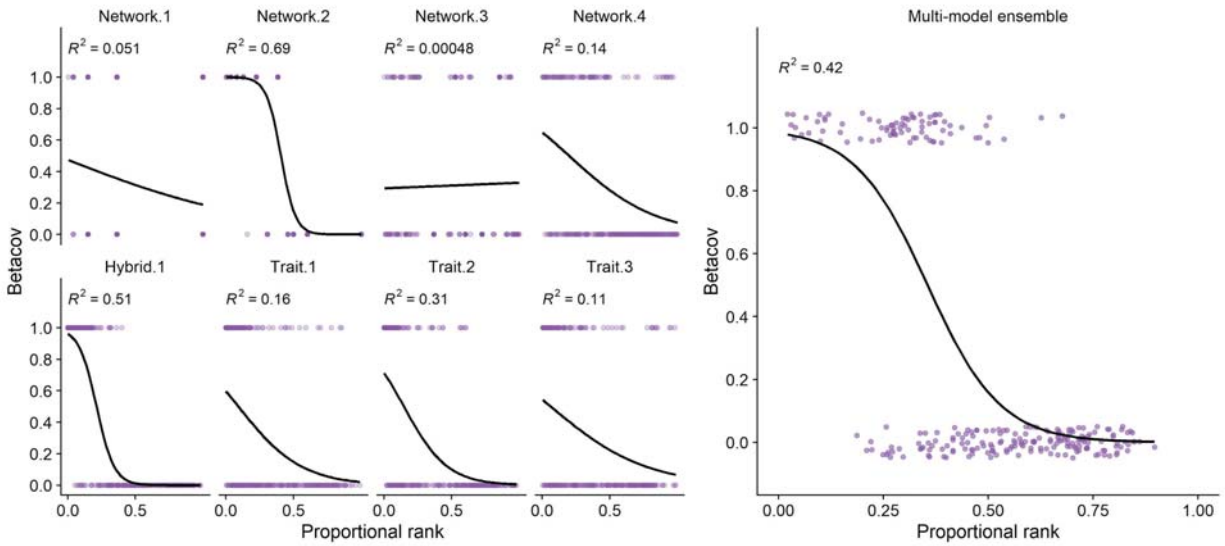
1012 **Supplemental Table 7. Model performance for hybrid model.** Area under receiver-operator curve for
1013 predicting interactions, or predicting for betacoronaviruses or for hosts, using bat traits, phylogeny or a
1014 combination of both.

1015

	interactions	Betacoronavirus	hosts (average)
traits	0.8784	0.6742	0.7304
phylogeny	0.8848	0.73389	0.7623
both	0.8975	0.7510	0.7909

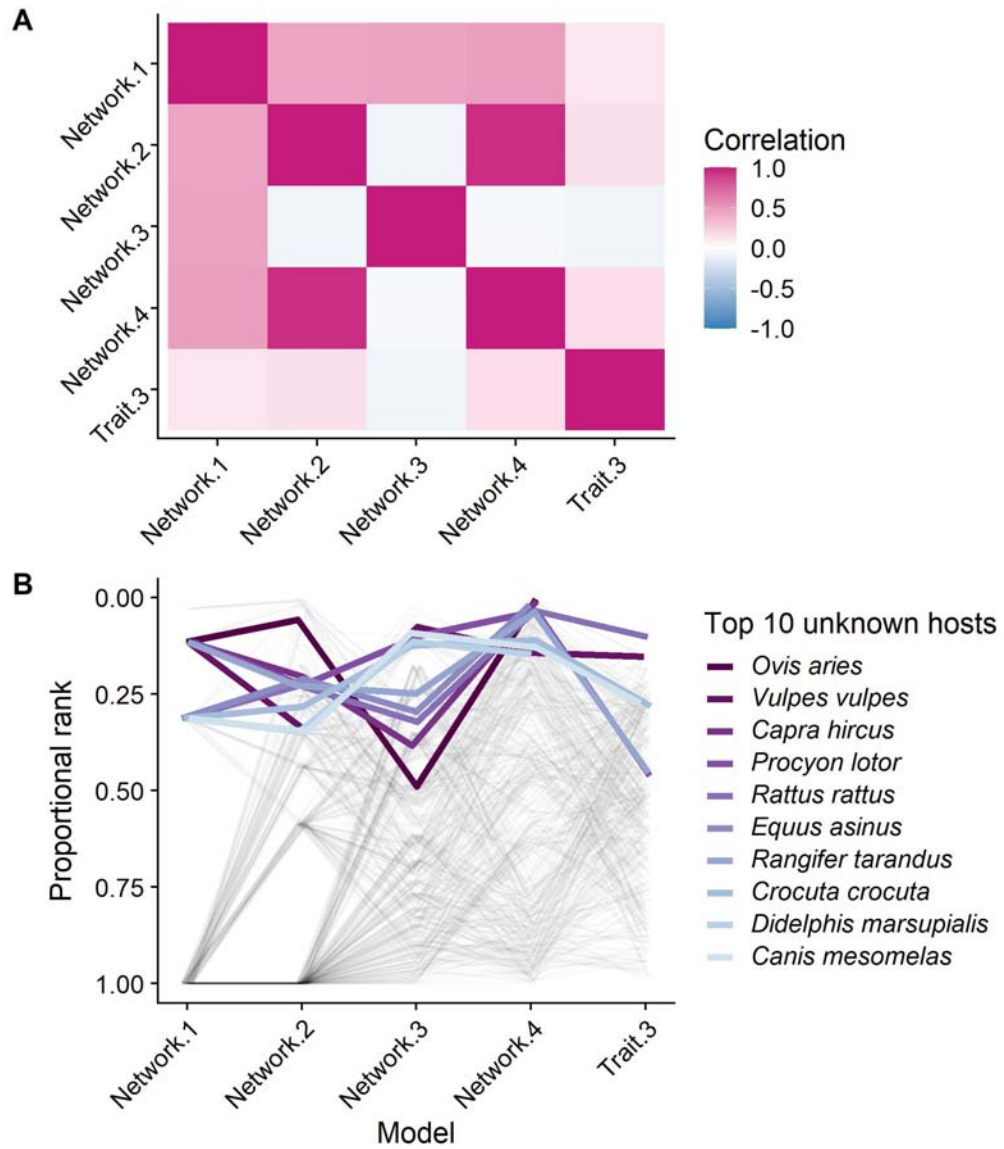
1016

1017 **Supplemental Figure 1.** Bat models perform more strongly together than in isolation. Curves show
1018 observed betacoronavirus hosts against predicted proportional ranks from eight individual models, and
1019 incorporated into one multi-model ensemble. Black lines show a binomial GLM fit to the predicted ranks
1020 against the recorded presence or absence of known betacoronavirus associations. Points are jittered to
1021 reduce overlap.
1022



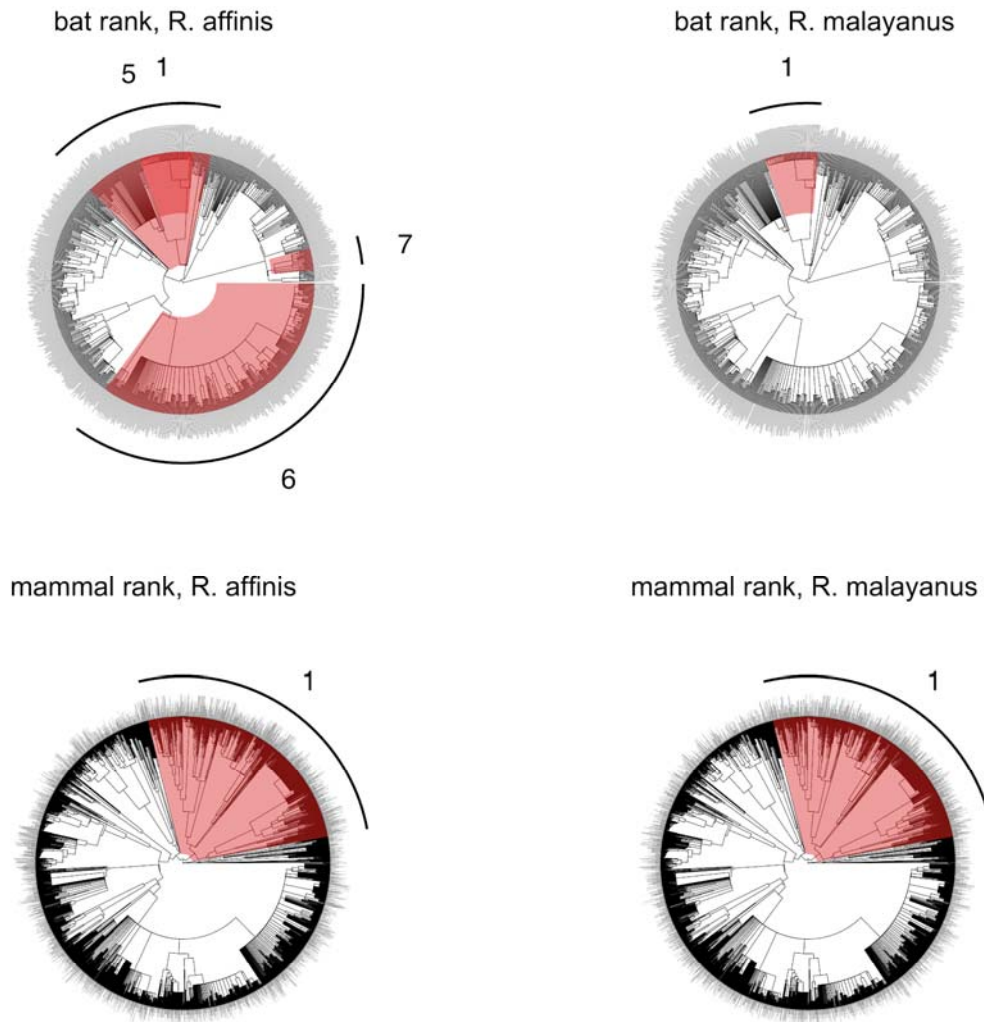
1023

1024 **Supplemental Figure 2.** Poor concordance among predictive models for mammal hosts of
1025 betacoronaviruses. The pairwise Spearman's rank correlations between models' ranked species-level
1026 predictions were generally low (A). In-sample predictions varied significantly and heavily prioritized
1027 domestic animals and well-studied hosts (B). The ten species with the highest mean proportional ranks
1028 across all models are highlighted in shades of purple. Only in-sample predictions are displayed because
1029 only one model (Trait-based 3) was able to predict out of sample for all mammals.
1030



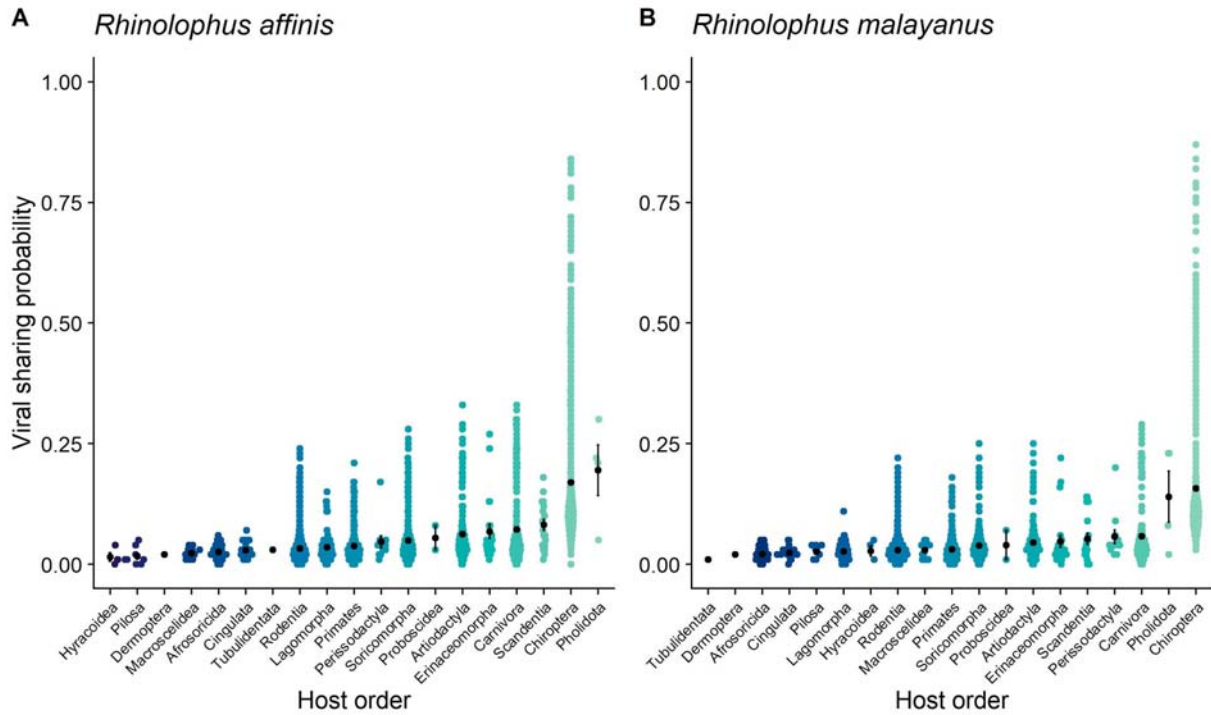
1031

1032 **Supplemental Figure 3.** Results of phylogenetic factorization applied to predicted ranks of virus sharing
1033 with *Rhinolophus affinis* and *Rhinolophus malayanus*. Colored regions indicate clades identified as
1034 significantly different in their predicted rank compared to the paraphyletic remainder; those more likely to
1035 share a virus with the *Rhinolophus* are shown in red, whereas those less likely to share a virus are shown
1036 in blue. Bar height indicates predicted rank (higher values = lower rank score, more likely share virus).
1037 Results are displayed for bats and remaining mammals separately. Mammal-wide clades with high
1038 propensities to share viruses with *R. affinis* based solely on their phylogeography included the treeshrews
1039 (Scandentia), Old World monkeys (Colobinae), and both tufted and barking deer (Muntiacini).
1040
1041



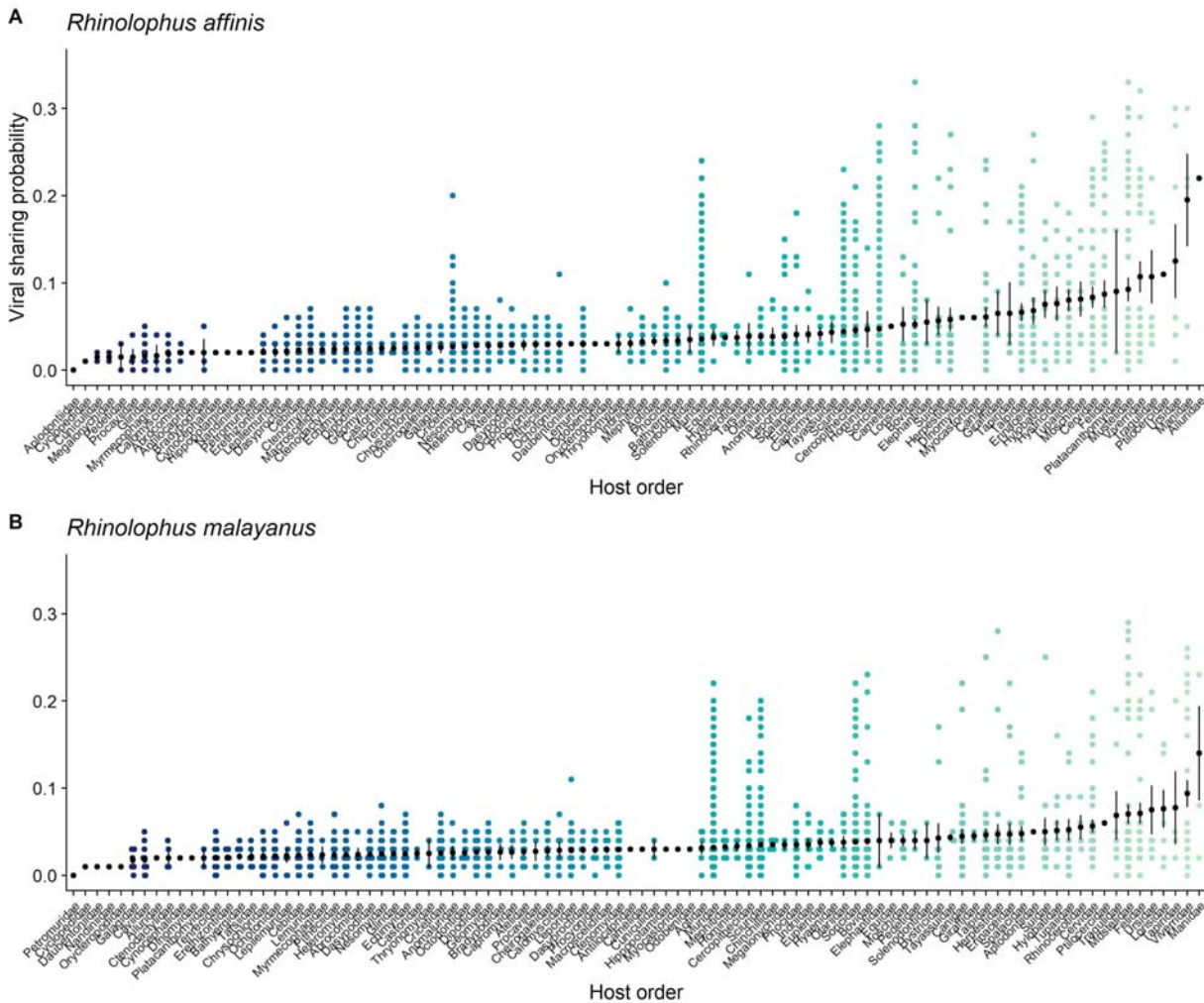
1042

1043 **Supplemental Figure 4.** Predicted species-level sharing probabilities of A) *Rhinolophus affinis* and B)
1044 *Rhinolophus malayanus*, calculated according to the phylogeographic viral sharing model⁴⁸. Each
1045 coloured point is a mammal species. Black points and error bars denote means and standard errors for
1046 each order. Mammal orders are arranged according to their mean sharing probability.
1047



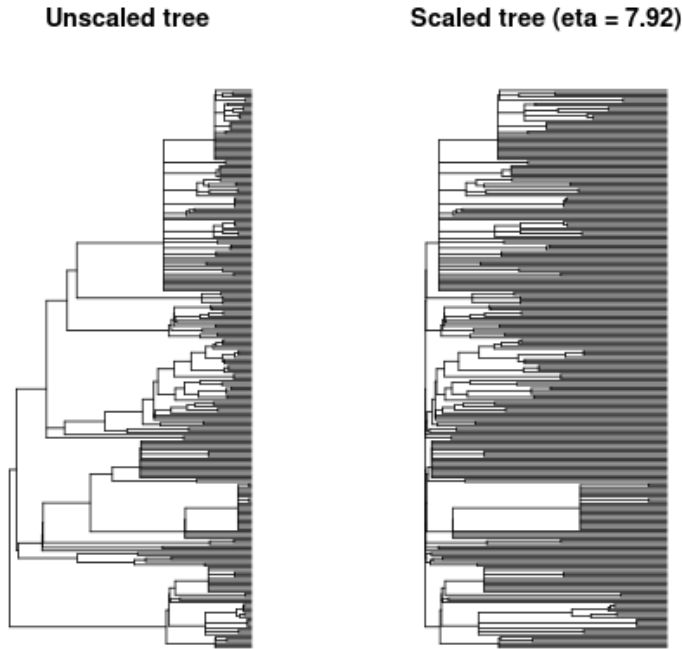
1048

1049 **Supplemental Figure 5.** Predicted species-level sharing probabilities of A) *Rhinolophus affinis* and B)
1050 *Rhinolophus malayanus*, calculated according to the phylogeographic viral sharing model[^]. Each
1051 coloured point is a mammal species. Black points and error bars denote means and standard errors for
1052 each order. Mammal orders are arranged according to their mean sharing probability.
1053



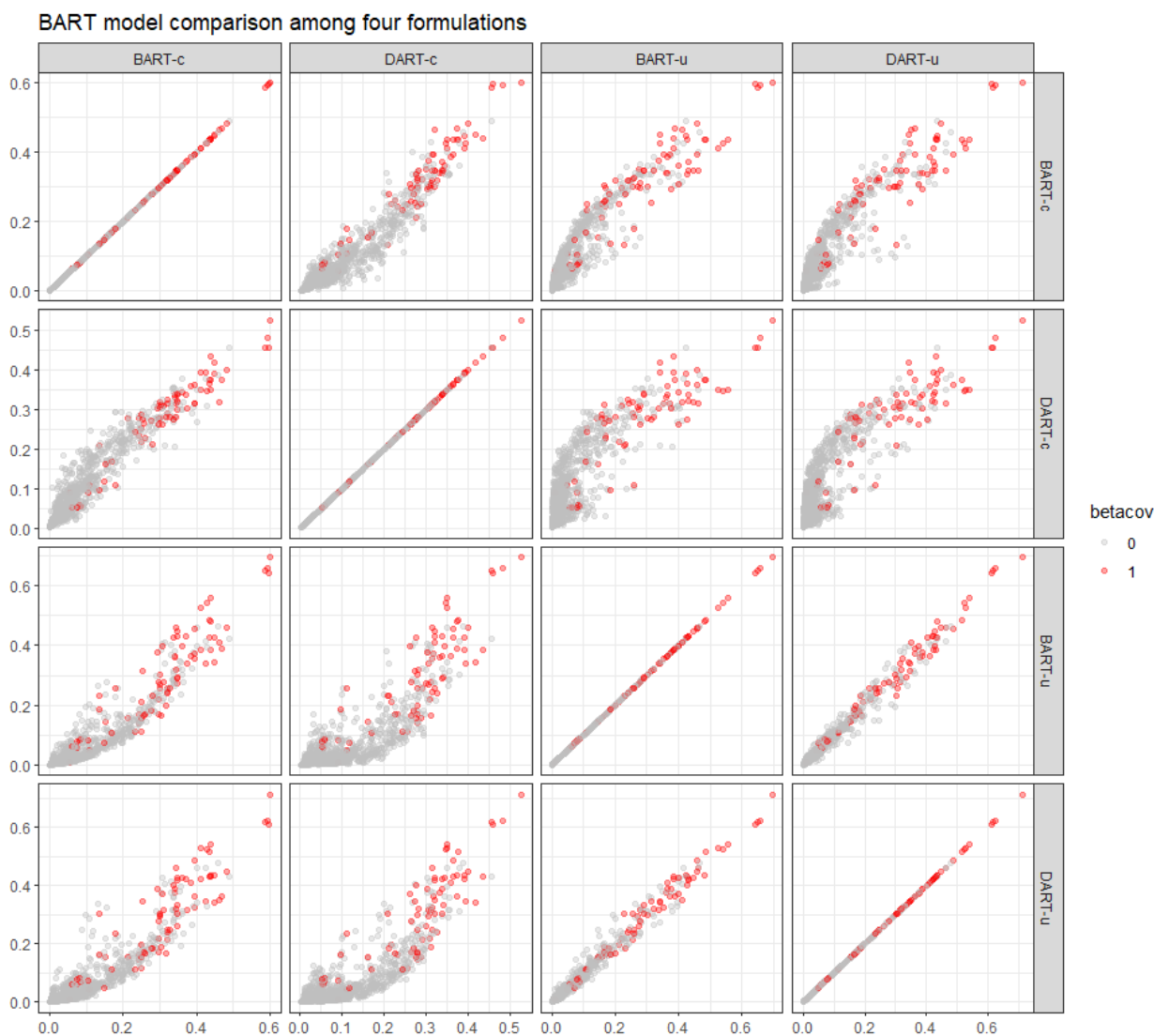
1054
1055

1056 **Supplemental Figure 6.** To account for uncertainty in the phylogenetic distances among hosts, the
1057 scaled-phylogeny model estimates a tree scaling parameter (η) based on an early-burst model of
1058 evolution. On the left is the unscaled bat phylogeny for the hosts in the bat-virus genera network, and on
1059 the right is the same tree rescaled according to mean estimated scaling parameter ($\eta = 7.92$). η values
1060 above 1 indicate accelerating evolution, suggesting less phylogenetic conservatism in host-virus
1061 associations among closely related taxa than would be predicted by a Brownian-motion model on the
1062 unscaled tree.
1063



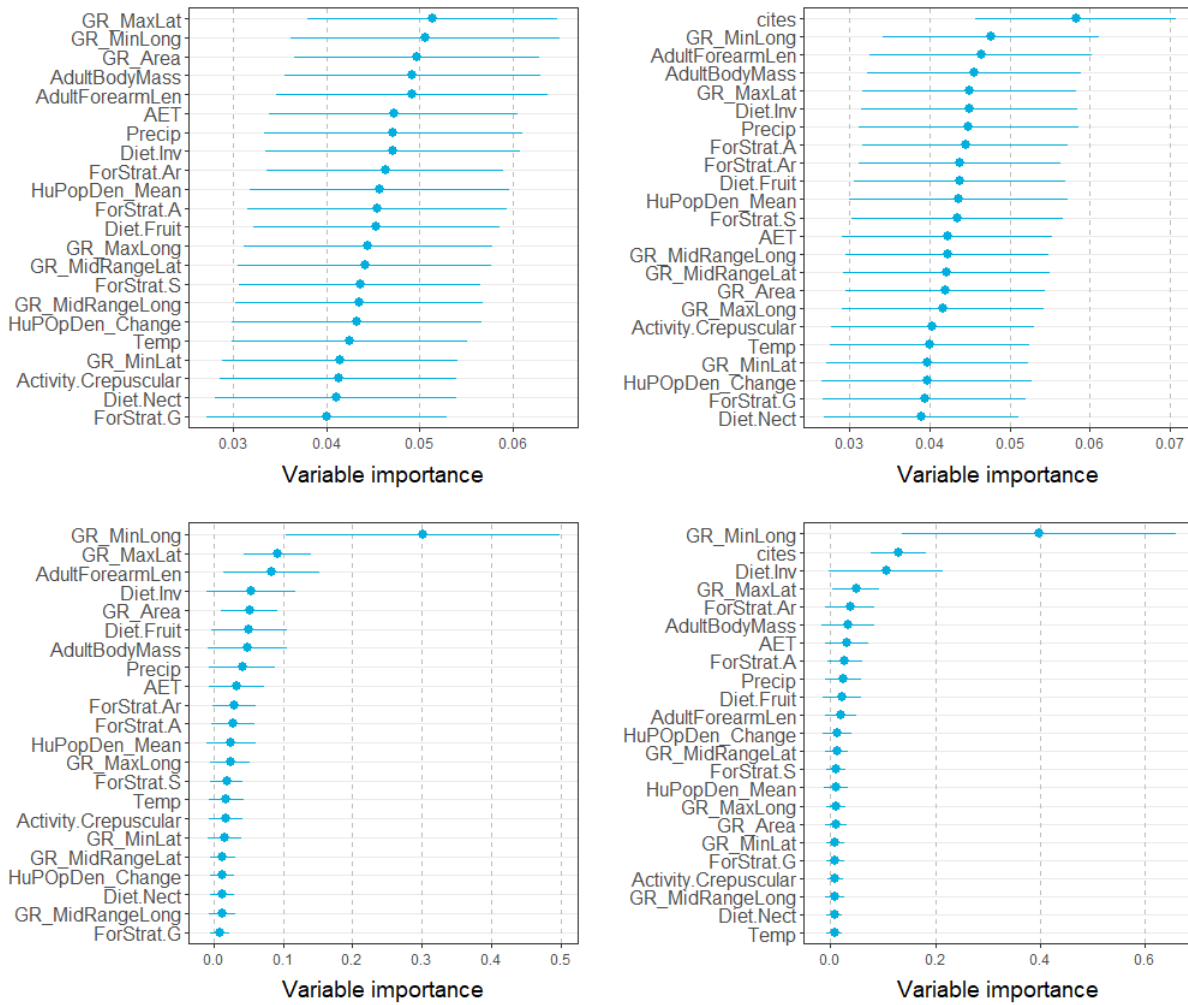
1064

1065 **Supplemental Figure 7.** Four formulations of Bayesian additive regression tree (BART) models produce
1066 slightly different results, but largely agree. Two models use baseline BART, while two models use a
1067 Dirichlet prior on variable importance (DART). Two are uncorrected for sampling bias (u) while two are
1068 corrected using citation counts (c). In the final main-text model ensemble, we present a DART model
1069 including correction for citation bias, which penalizes overfitting and spurious patterns two ways and
1070 leads to predictions with a lower total correlation with the data, but a still-high performance (AUC =
1071 0.90).
1072



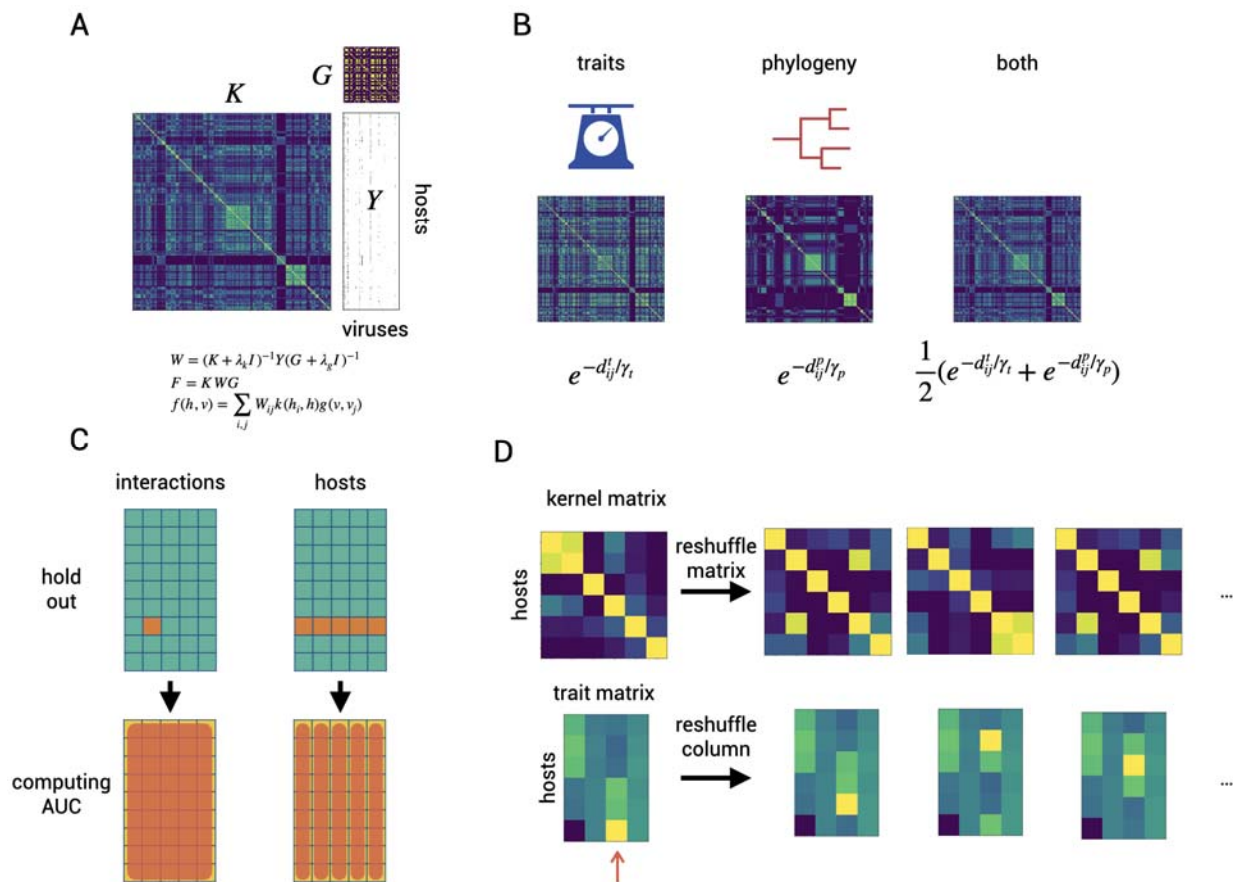
1073

1074 **Supplemental Figure 8.** Variable importance plots for the Bayesian additive regression tree models with
 1075 uniform variable importance prior (top) versus Dirichlet prior (bottom), without (left) and with (right)
 1076 correction for citations.



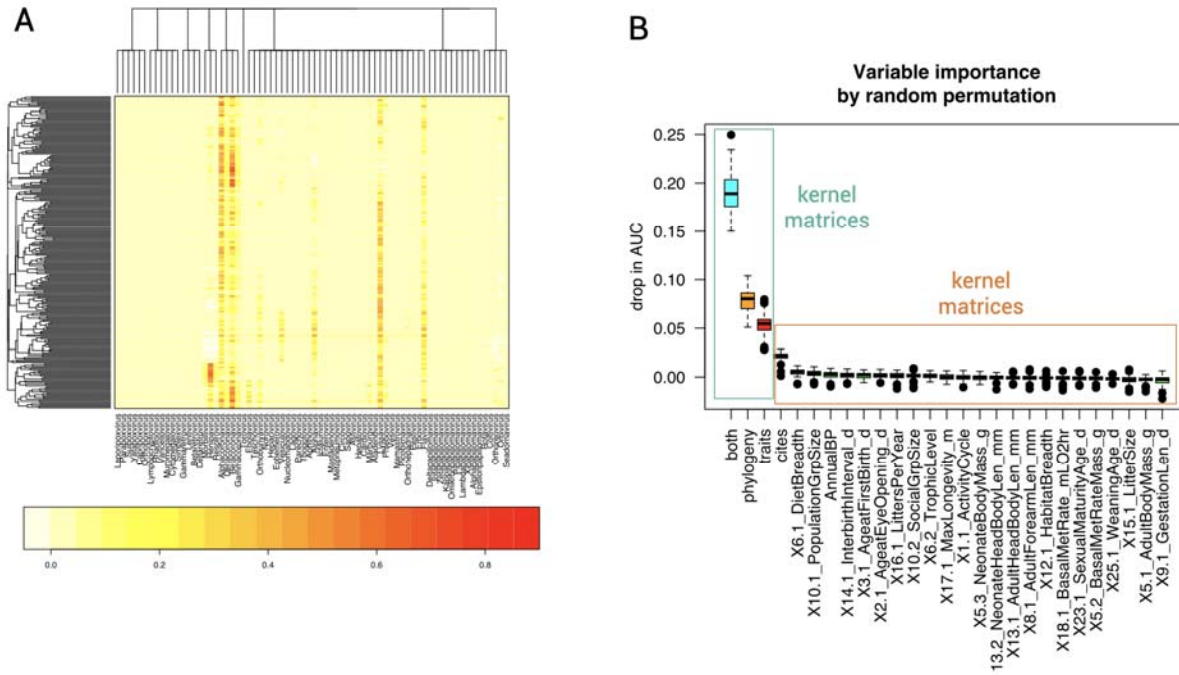
1077

1078 **Supplemental Figure 9.** Overview of the two-step kernel ridge regression (TSKRR) method. A. The
 1079 incidence matrix Y is modelled by two kernel matrices K and G , describing the hosts (here: bats) and
 1080 viruses, respectively. The equations for fitting the model and making predictions are given below. B.
 1081 Computing the host kernel matrix is done based on traits, phylogeny or both. A standard radial basis
 1082 kernel is used to transform a distance to a valid kernel. The two kernels are combined by averaging. C.
 1083 The models are validated by specialized cross-validation. The setting interactions leaves one interaction
 1084 out at a time and repeats this for every interaction. AUC is computed in micro-fashion. In the setting
 1085 hosts, we leave out one host at a time and compute the AUC per virus, averaging the individual results. D.
 1086 To assess the variable importance, we first randomly reshuffle the kernel matrices and monitor how the
 1087 performance of a model fitted on these data drops. In a more fine-grained approach, we just reshuffle a
 1088 single trait (indicated by the red arrow) and compare the resulting performances.
 1089



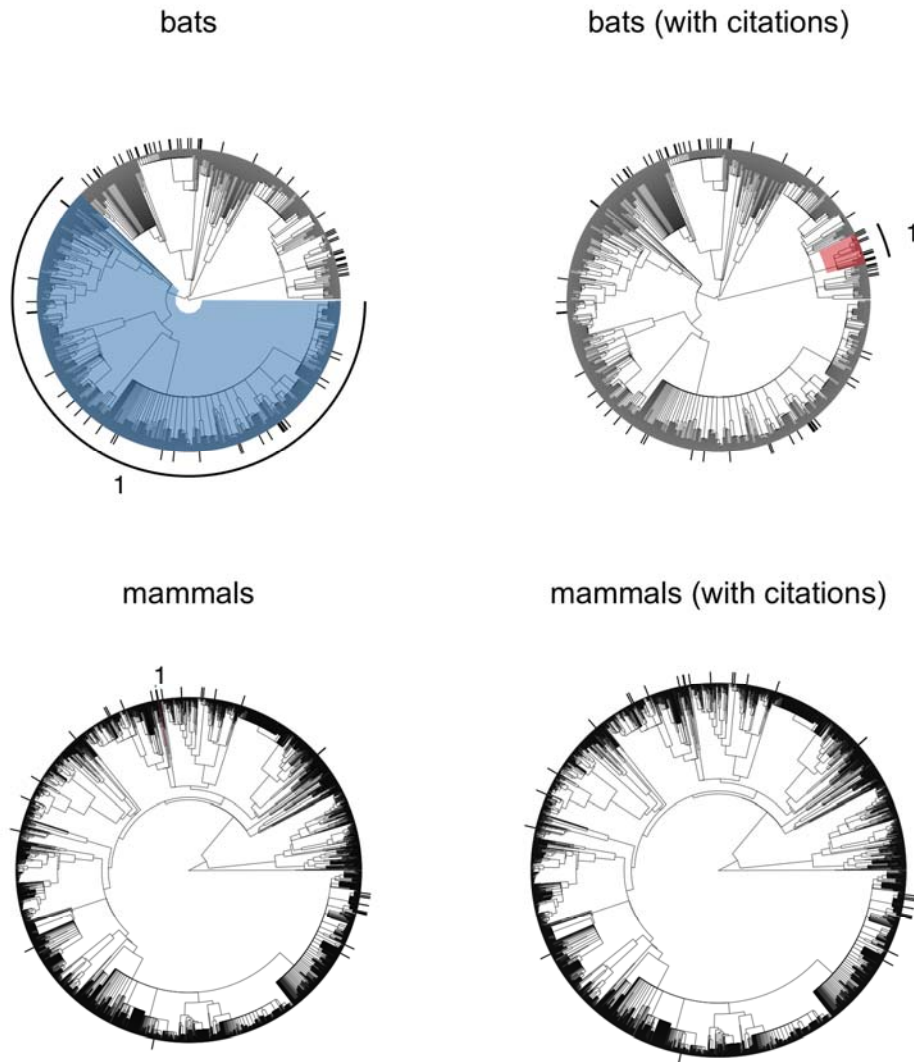
1090

1091 **Supplemental Figure 10.** Illustration of the TSKRR model and the variable importance scores. A.
1092 Heatmap of the prediction scores, together with a hierarchical clustering based on the kernel distance of K
1093 and G . B. Variable importance computed by random permutations of the kernel matrices or a single trait.



1094
1095
1096

1097 **Supplemental Figure 11.** Results of phylogenetic factorization applied to binomial betacoronavirus data
1098 across bats (top) and other mammals (bottom), using raw data (left) and after weighting by citation counts
1099 (right). Any significant clades (5% family-wise error rate) are displayed in colored shading on the
1100 phylogeny. Bars indicate betacoronavirus detection, and clades are colored by having more (red)
1101 (blue) positive species.
1102



1103

1104 **Bibliography**

- 1105 1. Anthony, S. J. *et al.* Global patterns in coronavirus diversity. *Virus Evol* **3**, vex012 (2017).
- 1106 2. Denison, M. R., Graham, R. L., Donaldson, E. F., Eckerle, L. D. & Baric, R. S. Coronaviruses: an
1107 RNA proofreading machine regulates replication fidelity and diversity. *RNA Biol.* **8**, 270–279
1108 (2011).
- 1109 3. Ren, W. *et al.* Full-length genome sequences of two SARS-like coronaviruses in horseshoe bats and
1110 genetic variation analysis. *J. Gen. Virol.* **87**, 3355–3359 (2006).
- 1111 4. Li, W. *et al.* Bats are natural reservoirs of SARS-like coronaviruses. *Science* **310**, 676–679 (2005).
- 1112 5. Yang, X.-L. *et al.* Isolation and Characterization of a Novel Bat Coronavirus Closely Related to the
1113 Direct Progenitor of Severe Acute Respiratory Syndrome Coronavirus. *J. Virol.* **90**, 3253–3256
1114 (2015).
- 1115 6. Zhou, P. *et al.* A pneumonia outbreak associated with a new coronavirus of probable bat origin.
1116 *Nature* **579**, 270–273 (2020).
- 1117 7. Guan, Y. *et al.* Isolation and characterization of viruses related to the SARS coronavirus from
1118 animals in southern China. *Science* **302**, 276–278 (2003).
- 1119 8. Hu, B. *et al.* Discovery of a rich gene pool of bat SARS-related coronaviruses provides new insights
1120 into the origin of SARS coronavirus. *PLoS Pathog.* **13**, e1006698 (2017).
- 1121 9. Memish, Z. A. *et al.* Middle East respiratory syndrome coronavirus in bats, Saudi Arabia. *Emerg.*
1122 *Infect. Dis.* **19**, 1819–1823 (2013).
- 1123 10. Wang, Q. *et al.* Bat origins of MERS-CoV supported by bat coronavirus HKU4 usage of human
1124 receptor CD26. *Cell Host Microbe* **16**, 328–337 (2014).
- 1125 11. Yang, Y. *et al.* Receptor usage and cell entry of bat coronavirus HKU4 provide insight into bat-to-
1126 human transmission of MERS coronavirus. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 12516–12521 (2014).
- 1127 12. Hu, B., Ge, X., Wang, L.-F. & Shi, Z. Bat origin of human coronaviruses. *Virology Journal* vol. 12
1128 (2015).
- 1129 13. Anthony, S. J. *et al.* Further Evidence for Bats as the Evolutionary Source of Middle East
1130 Respiratory Syndrome Coronavirus. *MBio* **8**, (2017).
- 1131 14. Anthony, S. J. *et al.* Coronaviruses in bats from Mexico. *J. Gen. Virol.* **94**, 1028–1038 (2013).
- 1132 15. Yang, L. *et al.* MERS-related betacoronavirus in *Vespertilio superans* bats, China. *Emerg. Infect.*
1133 *Dis.* **20**, 1260–1262 (2014).
- 1134 16. Zhou, H. *et al.* A novel bat coronavirus reveals natural insertions at the S1/S2 cleavage site of the
1135 Spike protein and a possible recombinant origin of HCoV-19. doi:10.1101/2020.03.02.974139.
- 1136 17. Nielsen, R., Wang, H. & Pipes, L. Synonymous mutations and the molecular evolution of SARS-
1137 Cov-2 origins. doi:10.1101/2020.04.20.052019.
- 1138 18. Lam, T. T.-Y. *et al.* Identifying SARS-CoV-2 related coronaviruses in Malayan pangolins. *Nature*
1139 (2020) doi:10.1038/s41586-020-2169-0.
- 1140 19. Xiao, K. *et al.* Isolation of SARS-CoV-2-related coronavirus from Malayan pangolins. *Nature*
1141 (2020) doi:10.1038/s41586-020-2313-x.
- 1142 20. Zhang, T., Wu, Q. & Zhang, Z. Probable Pangolin Origin of SARS-CoV-2 Associated with the
1143 COVID-19 Outbreak. *Curr. Biol.* **30**, 1578 (2020).
- 1144 21. Andersen, K. G., Rambaut, A., Ian Lipkin, W., Holmes, E. C. & Garry, R. F. The proximal origin of
1145 SARS-CoV-2. *Nature Medicine* vol. 26 450–452 (2020).
- 1146 22. Viana, M. *et al.* Assembling evidence for identifying reservoirs of infection. *Trends Ecol. Evol.* **29**,
1147 270–279 (2014).

- 1148 23. Plowright, R. K. *et al.* Prioritizing surveillance of Nipah virus in India. *PLoS Negl. Trop. Dis.* **13**,
1149 e0007393 (2019).
- 1150 24. Becker, D. J., Crowley, D. E., Washburne, A. D. & Plowright, R. K. Temporal and spatial
1151 limitations in global surveillance for bat filoviruses and henipaviruses. *Biol. Lett.* **15**, 20190423
1152 (2019).
- 1153 25. Becker, D. J., Washburne, A. D., Faust, C. L., Mordecai, E. A. & Plowright, R. K. The problem of
1154 scale in the prediction and management of pathogen spillover. *Philos. Trans. R. Soc. Lond. B Biol.*
1155 *Sci.* **374**, 20190224 (2019).
- 1156 26. Becker, D. J. & Han, B. A. The macroecology and evolution of avian competence for *Borrelia*
1157 *burgdorferi*. doi:10.1101/2020.04.15.040352.
- 1158 27. Han, B. A. *et al.* Undiscovered Bat Hosts of Filoviruses. *PLoS Negl. Trop. Dis.* **10**, e0004815
1159 (2016).
- 1160 28. Han, B. A. *et al.* Confronting data sparsity to identify potential sources of Zika virus spillover
1161 infection among primates. *Epidemics* **27**, 59–65 (2019).
- 1162 29. Washburne, A. D. *et al.* Taxonomic patterns in the zoonotic potential of mammalian viruses. *PeerJ*
1163 **6**, e5979 (2018).
- 1164 30. Olival, K. J. *et al.* Host and viral traits predict zoonotic spillover from mammals. *Nature* **546**, 646–
1165 650 (2017).
- 1166 31. Fritz, S. A., Bininda-Emonds, O. R. P. & Purvis, A. Geographical variation in predictors of
1167 mammalian extinction risk: big is bad, but only in the tropics. *Ecol. Lett.* **12**, 538–549 (2009).
- 1168 32. Jones, K. E. *et al.* PanTHERIA: a species-level database of life history, ecology, and geography of
1169 extant and recently extinct mammals. *Ecology* vol. 90 2648–2648 (2009).
- 1170 33. Wilman, H. *et al.* EltonTraits 1.0: Species-level foraging attributes of the world’s birds and
1171 mammals. *Ecology* vol. 95 2027–2027 (2014).
- 1172 34. Trifonova, N. *et al.* Spatio-temporal Bayesian network models with latent variables for revealing
1173 trophic dynamics and functional networks in fisheries ecology. *Ecological Informatics* vol. 30 142–
1174 158 (2015).
- 1175 35. Rohr, R. P., Scherer, H., Kehrl, P., Mazza, C. & Bersier, L.-F. Modeling food webs: exploring
1176 unexplained structure using latent traits. *Am. Nat.* **176**, 170–177 (2010).
- 1177 36. Dallas, T., Park, A. W. & Drake, J. M. Predicting cryptic links in host-parasite networks. *PLoS*
1178 *Comput. Biol.* **13**, e1005557 (2017).
- 1179 37. Han, B. A., Schmidt, J. P., Bowden, S. E. & Drake, J. M. Rodent reservoirs of future zoonotic
1180 diseases. *Proceedings of the National Academy of Sciences* vol. 112 7039–7044 (2015).
- 1181 38. Brandão, P. E. *et al.* A coronavirus detected in the vampire bat *Desmodus rotundus*. *Braz. J. Infect.*
1182 *Dis.* **12**, 466–468 (2008).
- 1183 39. Corman, V. M. *et al.* Highly diversified coronaviruses in neotropical bats. *J. Gen. Virol.* **94**, 1984–
1184 1994 (2013).
- 1185 40. Moreira-Soto, A. *et al.* Neotropical bats from Costa Rica harbour diverse coronaviruses. *Zoonoses*
1186 *Public Health* **62**, 501–505 (2015).
- 1187 41. Wang, L. *et al.* Discovery and genetic analysis of novel coronaviruses in least horseshoe bats in
1188 southwestern China. *Emerg. Microbes Infect.* **6**, e14 (2017).
- 1189 42. Lin, X.-D. *et al.* Extensive diversity of coronaviruses in bats from China. *Virology* vol. 507 1–10
1190 (2017).
- 1191 43. Wacharapluesadee, S. *et al.* Diversity of coronavirus in bats from Eastern Thailand. *Virol. J.* **12**, 57

- 1192 (2015).
- 1193 44. Guy, C., Ratcliffe, J. M. & Mideo, N. The influence of bat ecology on viral diversity and reservoir
1194 status. *Ecol. Evol.* **2008**, 209 (2020).
- 1195 45. Washburne, A. D. *et al.* Phylofactorization: a graph partitioning algorithm to identify phylogenetic
1196 scales of ecological data. *Ecol. Monogr.* **89**, e01353 (2019).
- 1197 46. Crowley, D., Becker, D., Washburne, A. & Plowright, R. Identifying Suspect Bat Reservoirs of
1198 Emerging Infections. *Vaccines* vol. 8 228 (2020).
- 1199 47. Damas, J., Hughes, G. M., Keough, K. C. & Painter, C. A. Broad Host Range of SARS-CoV-2
1200 Predicted by Comparative and Structural Analysis of ACE2 in Vertebrates. *bioRxiv* (2020).
- 1201 48. Albery, G. F., Eskew, E. A., Ross, N. & Olival, K. J. Predicting the global mammalian viral sharing
1202 network using phylogeography. *Nat. Commun.* **11**, 2260 (2020).
- 1203 49. Wang, M. *et al.* SARS-CoV Infection in a Restaurant from Palm Civet. *Emerging Infectious
1204 Diseases* vol. 11 1860–1865 (2005).
- 1205 50. Song, H.-D. *et al.* Cross-host evolution of severe acute respiratory syndrome coronavirus in palm
1206 civet and human. *Proc. Natl. Acad. Sci. U. S. A.* **102**, 2430–2435 (2005).
- 1207 51. Dougherty, E. R., Seidel, D. P., Carlson, C. J., Spiegel, O. & Getz, W. M. Going through the
1208 motions: incorporating movement analyses into disease research. *Ecol. Lett.* **21**, 588–604 (2018).
- 1209 52. Lehmann, D. *et al.* Pangolins and bats living together in underground burrows in Lopé National
1210 Park, Gabon. *African Journal of Ecology* (2020) doi:10.1111/aje.12759.
- 1211 53. Oreshkova, N. *et al.* SARS-CoV2 infection in farmed mink, Netherlands, April 2020.
1212 doi:10.1101/2020.05.18.101493.
- 1213 54. Kim, Y.-I. *et al.* Infection and Rapid Transmission of SARS-CoV-2 in Ferrets. *Cell Host Microbe
1214* **27**, 704–709.e2 (2020).
- 1215 55. Shi, J. *et al.* Susceptibility of ferrets, cats, dogs, and other domesticated animals to SARS–
1216 coronavirus 2. *Science* (2020) doi:10.1126/science.abb7015.
- 1217 56. Yang, X.-L. *et al.* Genetically Diverse Filoviruses in Rousettus and Eonycteris spp. Bats, China,
1218 2009 and 2015. *Emerg. Infect. Dis.* **23**, 482–486 (2017).
- 1219 57. Seifert, S. N. *et al.* Rousettus aegyptiacus Bats Do Not Support Productive Nipah Virus Replication.
1220 *The Journal of Infectious Diseases* vol. 221 S407–S413 (2020).
- 1221 58. Wacharapluesadee, S. *et al.* Longitudinal study of age-specific pattern of coronavirus infection in
1222 Lyle’s flying fox (*Pteropus lylei*) in Thailand. *Virol. J.* **15**, 38 (2018).
- 1223 59. Latinne, A. *et al.* Origin and cross-species transmission of bat coronaviruses in China.
1224 doi:10.1101/2020.05.31.116061.
- 1225 60. Yang, L. *et al.* Novel SARS-like betacoronaviruses in bats, China, 2011. *Emerg. Infect. Dis.* **19**,
1226 989–991 (2013).
- 1227 61. Geldenhuys, M. *et al.* A metagenomic viral discovery approach identifies potential zoonotic and
1228 novel mammalian viruses in Neoromicia bats within South Africa. *PLoS One* **13**, e0194527 (2018).
- 1229 62. Memish, Z. A. *et al.* Middle East respiratory syndrome coronavirus in bats, Saudi Arabia. *Emerg.
1230 Infect. Dis.* **19**, 1819–1823 (2013).
- 1231 63. Luo, Y. *et al.* Longitudinal Surveillance of Betacoronaviruses in Fruit Bats in Yunnan Province,
1232 China During 2009–2016. *Virol. Sin.* **33**, 87–95 (2018).
- 1233 64. Peel, A. J. *et al.* Synchronous shedding of multiple bat paramyxoviruses coincides with peak periods
1234 of Hendra virus spillover. *Emerg. Microbes Infect.* **8**, 1314–1323 (2019).
- 1235 65. de Souza Cortez J. L. Dunnum A. W. Ferguson F. A. Anwarali Khan D. L. Paul D. M. Reeder N. B.

- 1236 Simmons B. M. Thiers C. W. Thompson N S. Upham M. P. M. Vanhove P. W. Webala M. Weksler
1237 R. Yanagihara P. S. Soltis., C. J. A. S. A. B. A. J. B. C. A. C. B. M. B. Integrating biodiversity
1238 infrastructure into pathogen discovery and mitigation of epidemic infectious diseases. *Bioscience*
1239 (2020) doi:biaa064.
- 1240 66. Kingston, T. *et al.* Networking networks for global bat conservation. in *Bats in the Anthropocene:*
1241 *Conservation of Bats in a Changing World* 539–569 (Springer, Cham, 2016).
- 1242 67. Phelps, K. L. *et al.* Bat Research Networks and Viral Surveillance: Gaps and Opportunities in
1243 Western Asia. *Viruses* **11**, (2019).
- 1244 68. Teeling, E. C. *et al.* Bat Biology, Genomes, and the Bat1K Project: To Generate Chromosome-Level
1245 Genomes for All Living Bat Species. *Annu Rev Anim Biosci* **6**, 23–46 (2018).
- 1246 69. Mandl, J. N., Schneider, C., Schneider, D. S. & Baker, M. L. Going to Bat(s) for Studies of Disease
1247 Tolerance. *Front. Immunol.* **9**, 2112 (2018).
- 1248 70. Skirmuntt, E. C., Escalera-Zamudio, M., Teeling, E. C., Smith, A. & Katzourakis, A. The Potential
1249 Role of Endogenous Viral Elements in the Evolution of Bats as Reservoirs for Zoonotic Viruses.
1250 *Annu Rev Virol* (2020) doi:10.1146/annurev-virology-092818-015613.
- 1251 71. Jebb, D. *et al.* Six new reference-quality bat genomes illuminate the molecular basis and evolution of
1252 bat adaptations. (2019).
- 1253 72. Gervasi, S. S., Civitello, D. J., Kilvitis, H. J. & Martin, L. B. The context of host competence: a role
1254 for plasticity in host–parasite dynamics. *Trends Parasitol.* **31**, 419–425 (2015).
- 1255 73. Martin, L. B., Burgan, S. C., Adelman, J. S. & Gervasi, S. S. Host Competence: An Organismal Trait
1256 to Integrate Immunology and Epidemiology. *Integr. Comp. Biol.* **56**, 1225–1237 (2016).
- 1257 74. Callaway, E. & Cyranoski, D. Why snakes probably aren't spreading the new China virus. *Nature*
1258 (2020) doi:10.1038/d41586-020-00180-8.
- 1259 75. Gong, Y., Wen, G., Jiang, J. & Feng, X. Complete title: Codon bias analysis may be insufficient for
1260 identifying host(s) of a novel virus. *J. Med. Virol.* (2020) doi:10.1002/jmv.25977.
- 1261 76. Zhao, H. COVID-19 drives new threat to bats in China. *Science* **367**, 1436 (2020).
- 1262 77. Fenton, M. B. *et al.* Knowledge gaps about rabies transmission from vampire bats to humans. *Nature*
1263 *Ecology & Evolution* **4**, 517–518 (2020).
- 1264 78. López-Baucells, A., Rocha, R. & Fernández-Llamazares, Á. When bats go viral: negative framings
1265 in virological research imperil bat conservation. *Mamm. Rev.* **48**, 62–66 (2018).
- 1266 79. O'Shea, T. J., Cryan, P. M., Hayman, D. T. S., Plowright, R. K. & Streicker, D. G. Multiple
1267 mortality events in bats: a global review. *Mamm. Rev.* **46**, 175–190 (2016).
- 1268 80. MB Fenton, S Mubareka, SM Tsang, NB Simmons, DJ Becker. COVID-19 and threats to bats.
1269 *FACETS* in press (2020).
- 1270 81. Aguiar, L. M. S., Brito, D. & Machado, R. B. Do current vampire bat (*Desmodus rotundus*)
1271 population control practices pose a threat to Dekeyser's nectar bat's (*Lonchophylla dekeyseri*) long-
1272 term persistence in the Cerrado? *Acta Chiropt.* **12**, 275–282 (2010).
- 1273 82. Streicker, D. G. *et al.* Ecological and anthropogenic drivers of rabies exposure in vampire bats:
1274 implications for transmission and control. *Proc. Biol. Sci.* **279**, 3384–3392 (2012).
- 1275 83. Blackwood, J. C., Streicker, D. G., Altizer, S. & Rohani, P. Resolving the roles of immunity,
1276 pathogenesis, and immigration for rabies persistence in vampire bats. *Proc. Natl. Acad. Sci. U. S. A.*
1277 **110**, 20837–20842 (2013).
- 1278 84. Frick, W. F. *et al.* An emerging disease causes regional population collapse of a common North
1279 American bat species. *Science* **329**, 679–682 (2010).

- 1280 85. Sabir, J. S. M. *et al.* Co-circulation of three camel coronavirus species and recombination of MERS-
1281 CoVs in Saudi Arabia. *Science* **351**, 81–84 (2016).
- 1282 86. Guth, S., Visher, E., Boots, M. & Brook, C. E. Host phylogenetic distance drives trends in virus
1283 virulence and transmissibility across the animal–human interface. *Philos. Trans. R. Soc. Lond. B*
1284 *Biol. Sci.* **374**, 20190296 (2019).
- 1285 87. Redondo, R. A. F., Brina, L. P. S., Silva, R. F., Ditchfield, A. D. & Santos, F. R. Molecular
1286 systematics of the genus *Artibeus* (Chiroptera: Phyllostomidae). *Mol. Phylogenet. Evol.* **49**, 44–58
1287 (2008).
- 1288 88. Bouchard, S. *Chaerephon pumilus*. *Mammalian Species* 1–6 (1998).
- 1289 89. Hooper, S. R., Van Den Bussche, R. A. & Horáček, I. Generic Status of the American Pipistrelles
1290 (Vespertilionidae) with Description of a New Genus. *J. Mammal.* **87**, 981–992 (2006).
- 1291 90. Desjardins-Proulx, P., Laigle, I., Poisot, T. & Gravel, D. Ecological interactions and the Netflix
1292 problem. *PeerJ* **5**, e3644 (2017).
- 1293 91. Stock, M., Poisot, T., Waegeman, W. & De Baets, B. Linear filtering reveals false negatives in
1294 species interaction data. *Sci. Rep.* **7**, 45908 (2017).
- 1295 92. Drake, J. M. & Richards, R. L. Estimating environmental suitability. *Ecosphere* vol. 9 e02373
1296 (2018).
- 1297 93. Dallas, T. A., Carlson, C. J. & Poisot, T. Testing predictability of disease outbreaks with a simple
1298 model of pathogen biogeography. *R Soc Open Sci* **6**, 190883 (2019).
- 1299 94. Elmasri, M., Farrell, M. J., Jonathan Davies, T. & Stephens, D. A. A hierarchical Bayesian model for
1300 predicting ecological interactions using scaled evolutionary relationships. *The Annals of Applied*
1301 *Statistics* vol. 14 221–240 (2020).
- 1302 95. Cadotte, M. W. *et al.* Phylogenetic diversity metrics for ecological communities: integrating species
1303 richness, abundance and evolutionary history. *Ecol. Lett.* **13**, 96–105 (2010).
- 1304 96. Park, A. W. *et al.* Characterizing the phylogenetic specialism–generalism spectrum of mammal
1305 parasites. *Proceedings of the Royal Society B: Biological Sciences* vol. 285 20172613 (2018).
- 1306 97. Harvey, P. H. & Pagel, M. D. *The comparative method in evolutionary biology*. (Oxford University
1307 Press, USA, 1998).
- 1308 98. Harmon, L. J. *et al.* Early bursts of body size and shape evolution are rare in comparative data.
1309 *Evolution* **64**, 2385–2396 (2010).
- 1310 99. Mollentze, N. & Streicker, D. G. Viral zoonotic risk is homogenous among taxonomic orders of
1311 mammalian and avian reservoir hosts. *Proc. Natl. Acad. Sci. U. S. A.* **117**, 9423–9430 (2020).
- 1312 100. Schmidt, J. P. *et al.* Ecological indicators of mammal exposure to Ebolavirus. *Philos. Trans. R. Soc.*
1313 *Lond. B Biol. Sci.* **374**, 20180337 (2019).
- 1314 101. Pandit, P. S. *et al.* Predicting wildlife reservoirs and global vulnerability to zoonotic Flaviviruses.
1315 *Nat. Commun.* **9**, 5425 (2018).
- 1316 102. Evans, M. V., Dallas, T. A., Han, B. A., Murdock, C. C. & Drake, J. M. Data-driven identification of
1317 potential Zika virus vectors. *Elife* **6**, (2017).
- 1318 103. Yang, L. H. & Han, B. A. Data-driven predictions and novel hypotheses about zoonotic tick vectors
1319 from the genus *Ixodes*. *BMC Ecol.* **18**, 7 (2018).
- 1320 104. Elith, J., Leathwick, J. R. & Hastie, T. A working guide to boosted regression trees. *J. Anim. Ecol.*
1321 **77**, 802–813 (2008).
- 1322 105. Carlson, C. J. *embarcadero*: Species distribution modelling with Bayesian additive regression trees in
1323 R. doi:10.1101/774604.

- 1324 106. Chipman, H. A., George, E. I. & McCulloch, R. E. BART: Bayesian additive regression trees. *The*
1325 *Annals of Applied Statistics* vol. 4 266–298 (2010).
- 1326 107. Chen, W. *et al.* The illegal exploitation of hog badgers (*Arctonyx collaris*) in China: genetic
1327 evidence exposes regional population impacts. *Conservation Genetics Resources* vol. 7 697–704
1328 (2015).
- 1329 108. Stock, M., Pahikkala, T., Airola, A., De Baets, B. & Waegeman, W. A Comparative Study of
1330 Pairwise Learning Methods Based on Kernel Ridge Regression. *Neural Comput.* **30**, 2245–2283
1331 (2018).
- 1332 109. Stock, M., Pahikkala, T., Airola, A., Waegeman, W. & De Baets, B. Algebraic shortcuts for leave-
1333 one-out cross-validation in supervised network inference. *Brief. Bioinform.* (2018)
1334 doi:10.1093/bib/bby095.
- 1335 110. Stekhoven, D. J. & Bühlmann, P. MissForest—non-parametric missing value imputation for mixed-
1336 type data. *Bioinformatics* **28**, 112–118 (2012).
- 1337 111. Schölkopf, B., Smola, A. J., Managing Director of the Max Planck Institute for Biological
1338 Cybernetics in Tübingen Germany Profe Bernhard Scholkopf & Bach, F. *Learning with Kernels:*
1339 *Support Vector Machines, Regularization, Optimization, and Beyond.* (MIT Press, 2002).
- 1340 112. Foley, N. M., Goodman, S. M., Whelan, C. V., Puechmaille, S. J. & Teeling, E. Towards navigating
1341 the Minotaur’s labyrinth: cryptic diversity and taxonomic revision within the speciose genus
1342 *Hipposideros* (Hipposideridae). *Acta Chiropt.* **19**, 1–18 (2017).