# Predicting the global mammalian viral sharing network using phylogeography

Gregory F Albery[1,2,3*], Evan A Eskew[1], Noam Ross[1], Kevin J Olival[1*]

1. EcoHealth Alliance, New York, NY, USA
2. Institute of Evolutionary Biology, University of Edinburgh, Edinburgh, Scotland
3. Department of Biology, Georgetown University, Washington, DC, USA

*Corresponding authors: gfalbery@gmail.com; olival@ecohealthalliance.org

# Abstract

Understanding interspecific viral transmission is key to understanding viral ecology and evolution, disease spillover into humans, and the consequences of global change. Prior work has demonstrated that macroecological factors drive viral sharing in some mammalian groups, but analyses have never attempted to predict viral sharing in a pan-mammalian context. Here we show that host phylogenetic similarity and geographic range overlap are strong, nonlinear predictors of viral sharing among species across the entire mammal class. Using these traits, we predict global viral sharing patterns across 4196 mammal species and show that our simulated network successfully predicts viral sharing and reservoir host status using internal validation and an external dataset. We predict high rates of mammalian viral sharing in the tropics, particularly among rodents and bats, and that within- and between-order sharing differs geographically and taxonomically. Our results emphasize the importance of macroecological factors in shaping mammalian viral communities, and provide a robust, general model to predict viral host range and guide pathogen surveillance and conservation efforts.

27　Most emerging human viruses originate in wild mammals, so understanding the drivers of

28　interspecific viral transmission in these taxa is an important public health research priority[1,2].

29　Despite a rapidly expanding knowledge base, the mammalian viruses known to science

30　remain taxonomically biased and limited in scope, likely comprising less than 1% of the

31　complete mammalian virome[3,4]. Furthermore, host range is inadequately characterized even

32　for the best-studied viruses[5–7]. To help prioritise viral discovery efforts and zoonotic disease

33　surveillance in wildlife, studies have revealed high (zoonotic) parasite diversity in certain

34　host taxa, such as rodents and bats[5,8], and/or linked parasite diversity with host phenotypic

35　traits such as reproductive output[9,10]. Viral diversity has also been associated with host

36　macroecological traits, including geographic range size[11] and sympatry with other mammals[5].

37　The rationale for investigating viral diversity is that species with more viruses will generate

38　more opportunities for viral transmission to other species, including humans. However, in

39　order to infect a new host species, a virus must transmit, invade, and potentially replicate

40　within the novel host[12]. Each of these processes becomes less likely if the two hosts differ

41　more in terms of their geographic range, behaviour, and/or biochemistry (i.e., cellular

42　receptors allowing viral attachment and invasion)[12,13]. Consequently, the probability that a

43　pair of hosts will share a virus is shaped both by the species' underlying viral diversity and by

44　species interactions represented by pairwise measures such as spatial overlap, phylogenetic

45　relatedness, and ecological similarity[14–16].

46

47　Previous investigations into pairwise determinants of viral sharing have been limited to one

48　or two host orders (e.g., bats[17,18], primates[19], ungulates[16], and carnivores[14,16]), while

49　sometimes lumping together different types of pathogen (e.g., helminths, viruses, and

50　bacteria). Viruses are sometimes shared across large host phylogenetic distances (e.g., Nipah

51　virus in bats and pigs, among many others[20,21]), requiring a broader understanding of viral

52　sharing across mammals to predict patterns at different taxonomic and geographic scales. In

53　addition, many mammalian orders have yet to be investigated in these analyses – most

54　notably rodents, which are highly diverse and host important zoonotic viruses[5,8]. In addition,

55　although phylogenetic and geographic viral sharing effects have been empirically

56　demonstrated, the models have not yet been applied to validate viral sharing predictions using

57　external datasets or make inferences about mammals with no known viral associations. If

58　geographic and phylogenetic effects on viral sharing are as ubiquitous as they seem, these

59　variables alone could provide a useful baseline model of viral sharing applicable across the
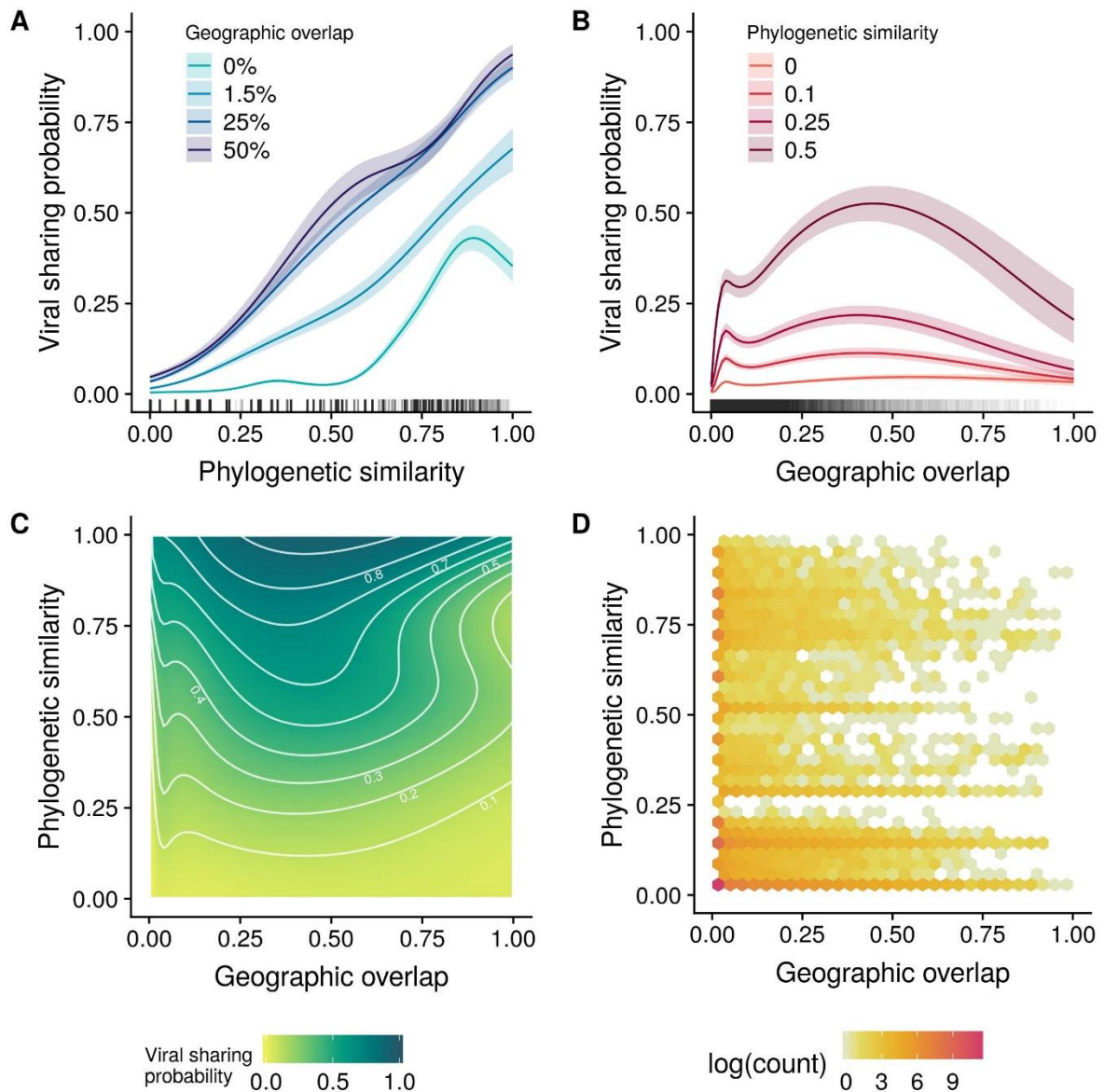
60　mammal class.

61 Here, we analyse pairwise viral sharing using a novel, conservative modelling approach

62 designed to partition the contribution of species-level traits from pairwise phylogeographic

63 traits. This method of analysis stands in contrast to previous studies of mammalian viral

64 sharing which have mainly focussed on host-level traits, and importantly buffers against

65 certain inherent biases in the observed viral sharing network, including host sampling bias,

66 when making predictions.

# Results and Discussion

## Predictors of viral sharing

69 We fitted a model designed to partition the contribution of species-level effects and pairwise

70 similarity measures to mammalian viral sharing probability. We used a published database of

71 1920 mammal-virus associations (excluding humans) as a training dataset[5]. These data

72 included 591 wild mammal species, equalling 174345 pairwise host species combinations,

73 with 6.4% connectance – that is, 6.4% of species pairs shared at least one virus. We used a

74 generalised additive mixed model (GAMM) framework, including a species-level effect in

75 our model as a multi-membership random effect, capturing variation in each species'

76 connectedness and underlying viral diversity (see Methods). Overall, our model accounted

77 for 44.8% of the total deviance in pairwise viral sharing, with 51.1% of this explained

78 deviance attributable to the identities of the species involved (i.e., the species-level effect).

79 Our model structure was effective at controlling for species-level variation in our dataset: i.e.,

80 the term had a strong impact on the centrality of each species when we simulated networks

81 using just these parameters (Figure SI1). This observation suggests that ~50% of the dyadic

82 structure of observed viral sharing networks (in contrast to the true underlying network) is

83 determined by uneven sampling and concentration on specific species, and the remainder by

84 macroecological processes.

85

Figure 1: Viral sharing GAMM model outputs and data distribution. A: predicted viral sharing probability increases with increasing phylogenetic relatedness; the different coloured lines represent different geographic overlap values. B: predicted viral sharing probability increases with increasing geographic overlap; the different coloured lines represent different phylogenetic relatedness values. C: the geographic overlap:phylogenetic similarity interaction surface, where the darker colours represent increased probability of viral sharing. White contour lines denote 10% increments of sharing probability. Labels have been removed from some contours to avoid overplotting. D: hexagonal bin chart displaying the data distribution, which was highly aggregated at low values of phylogenetic similarity and especially of geographic overlap.

As expected, increasing host phylogenetic similarity and geographic overlap were associated with increased probability of viral sharing across mammals, together accounting for the

99    remaining 49% of explained model deviance (Figure 1A-C). Geography, phylogeny, and their

100    interaction all showed strong nonlinear effects, with geographic overlap in particular driving

101    a rapid increase in viral sharing that began at ~0-5% range overlap values, peaked at 50%

102    overlap values, and then levelled off (Figure 1B). This effect closely mirrors previous

103    observations of strong, nonlinear effects of geographic and phylogenetic similarity

104    determining within-order viral sharing[14,16–19]. Although occupying little of the visual space

105    within the model presentation, 93% of mammal pairs had less than 5% spatial overlap (Figure

106    1B,D). The great majority (86%) of mammal pairs in our dataset did not overlap

107    geographically and rarely shared viruses unless phylogenetic similarity exceeded ~0.5

108    (Figure 1A). This phylogenetic distance corresponds roughly to order-level similarity; that is,

109    if two species did not overlap in space, it was highly unlikely that they shared a virus unless

110    they were within the same taxonomic order (8% of pairs). Notably, phylogenetic similarity

111    accounted for more than twice as much model deviance as did spatial overlap (33.8% vs

112    14.4%). The greater importance of phylogeny relative to geography contrasts with previous

113    analyses concerning viral sharing in primates[19] and ungulates[16], likely reflecting the wider

114    phylogenetic range of hosts considered here. This finding supports the important role of

115    mammalian evolutionary history in shaping contemporary patterns of viral sharing and

116    diversity[5,22].

117

118    In contrast to geography and phylogeny, minimum citation count and domestication status

119    accounted for a vanishingly small amount of the deviance in viral sharing probability (0.2%

120    and 0.1%, respectively) even though they have important effects on observed viral diversity

121    in this dataset[5]. Their impacts on viral sharing may have been largely accounted for by

122    species-level random effects.

123

124    Our use of a pan-mammalian viral sharing dataset with a large sample size allowed us to

125    investigate how geographic overlap and phylogenetic similarity affect viral sharing across

126    different viral subgroups. These subgroups included RNA viruses, vector-borne RNA viruses,

127    non-vector-borne RNA viruses, and DNA viruses. The importance of geographic overlap

128    varied widely across all groups of viruses (Figure SI2; Table SI1), while the influence of host

129    phylogenetic relatedness was more consistent (Figure SI3; Table SI1). Generally, host

130    phylogeny was more important in determining sharing of DNA viruses than it was for RNA

131    viruses, while space sharing was more important for vector-borne RNA viruses, and less so

132    for non-vector-borne RNA viruses. These results likely reflect important aspects of viral

133    ecology, transmission, and evolution: for example, RNA viruses are fast-evolving, allowing

134    them to more quickly adapt to novel hosts, such that phylogenetic distances are less important

135    in determining viral sharing patterns[23]. Conversely, DNA viruses are more evolutionarily

136    constrained, with an evolutionary rate typically <1% that of RNA viruses, such that

137    phylogenetic distance between hosts presents a more significant obstacle for sharing of DNA

138    viruses[24]. The profound importance of geographic overlap in shaping the viral sharing

139    network for vector-borne RNA viruses (Figure SI3) likely emerges from the geographic

140    distributions and ecological constraints placed on vectors, lending further support to efforts to

141    model the global spread of arboviruses by predicting changes in their vectors' distributions

142    and ecological niches[25,26]. Generally, the fact that viral sharing across different viral

143    subgroups was predicted by different macroecological relationships suggests they should be

144    examined separately in future analyses where possible.


## Predicting pan-mammalian viral sharing

146    Previous trait-based approaches to predict viral sharing and reservoir hosts have been

147    hindered by incomplete and inconsistent characterization of traits central to those modelling

148    efforts. In contrast, spatial distributions and phylogenetic data are readily available and

149    uniformly quantified for the vast majority of mammals and, as we have shown, are reliable

150    predictors of viral sharing (>20% of total deviance). Thus, we used our GAMM estimates to

151    predict unobserved global viral sharing patterns across 8.8 million mammal-mammal pairs

152    using a database of geographic distributions[27] and a recent mammalian supertree[28] (see

153    Methods). The predicted network included 4196 (non-human) Eutherian mammals with

154    available data, 591 of which were recorded with viral associations in our training data. We

155    calculated each species' predicted degree centrality, as a simple and interpretable network-

156    derived measure of viral sharing: that is, the number of other mammal species a given

157    mammal species is expected to share at least one virus with. We identified geographic and

158    taxonomic trends in degree centrality, validated our predicted sharing network using an

159    external dataset, and simulated reservoir identification to assess host predictability for focal

160    viruses (see Methods).

161

162    We confirmed that our modelled network recapitulated expected patterns of viral sharing

163    using the Enhanced Infectious Diseases Database (EID2) as an external dataset[29]. This dataset

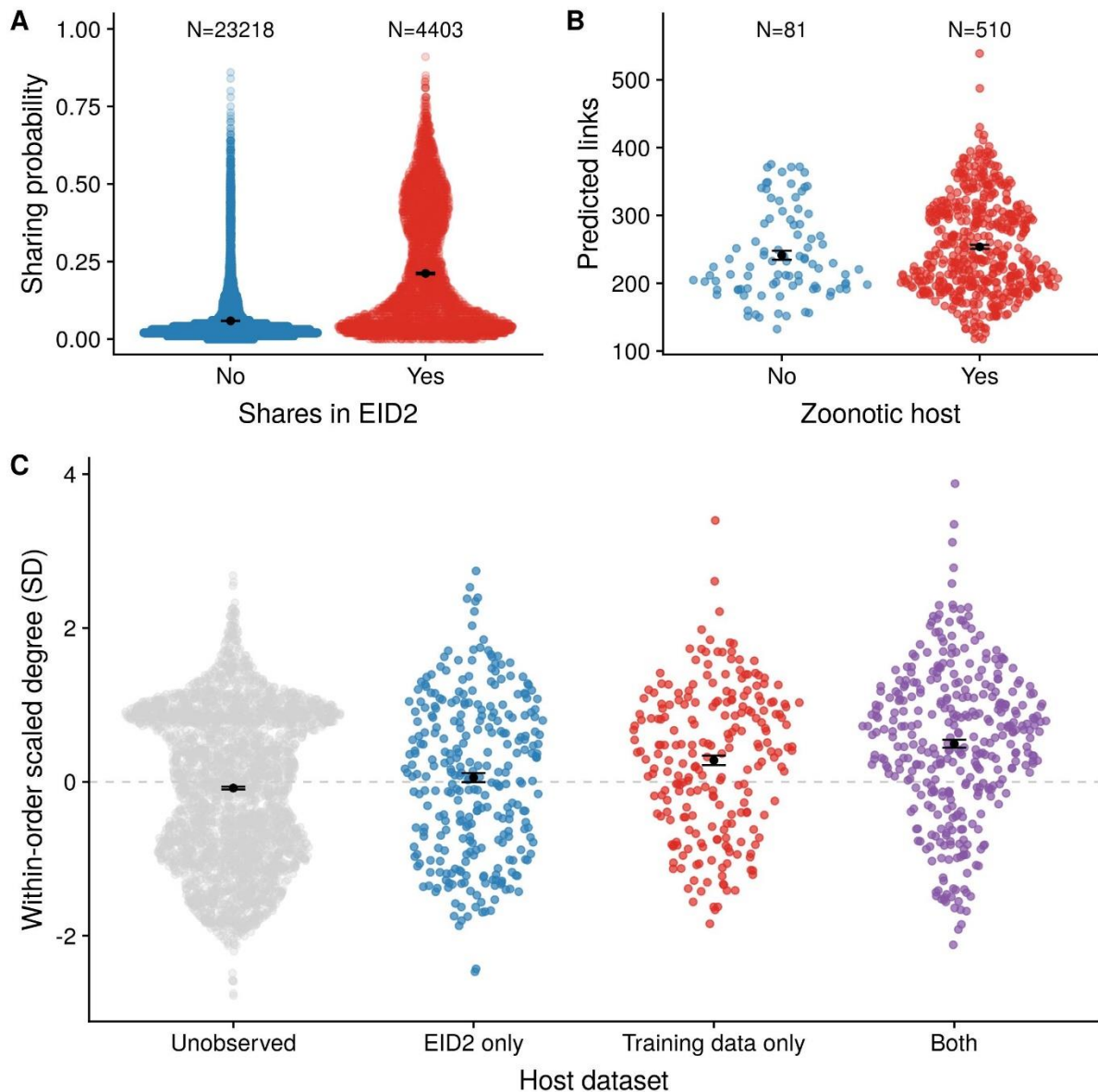164    was constructed by mining web-based sequence data to identify host-pathogen associations,

165    many of which are mammal-virus interactions[29]. Pairs of species that share viruses in EID2,

166    but which were not in our training dataset (see Methods), had a much higher mean sharing

167    probability in our predicted network (20% versus 5%; Figure 2A). In addition, more central

168    species in the predicted network were more likely to have been observed with a virus,

169    whether zoonotic (Figure 2B) or non-zoonotic (Figure 2C), implying that the predicted

170    network accurately captured realised potential for viral sharing and zoonotic spillover. This

171    finding concurs with similar work in primates which demonstrated that high centrality in

172    primate-parasite networks is associated with carriage of zoonoses[30]. We corroborate these

173    findings considering all mammal-mammal viral sharing links, not just zoonotic links, and

174    show that for each mammalian order, species with higher degree centrality in our predicted

175    network are more likely to have been observed with viruses in the EID2 dataset (Figure 2C;

176    Figure SI4). It is possible that species with higher centrality in the global viral sharing

177    network are more important for viral sharing, and thus have been more likely to be observed

178    with a (zoonotic) virus. Species that are more central in our predicted network could therefore

179    be prioritised for zoonotic surveillance or sampling in the event of viral outbreaks with

180    unknown mammalian origins. Given that mammal diversity predicts patterns of livestock

181    disease[31] and zoonoses[32], the geographic patterns of degree centrality predicted here (Figure

182    3, Figure SI5; see below) could also be used as a coarse predictor of viral disease risk to

183    livestock and human health, providing additional insights that emerge from the joint,

184    nonlinear effects of geography and phylogeny as opposed to examination of their effects in

185    isolation. Similarly, where there is limited knowledge of mammalian host range for newly-

186    discovered viruses, our modelled network can be used to prioritise the sampling of additional

187    species for viral surveillance.

188

189    The high predicted centrality of known hosts may be due partly to selective sampling (i.e.,

190    viral researchers are more likely to sample wide-ranging and common host species that also

191    share viruses with many other species[10,20]). This possibility is supported by the increased

192    degree centrality for species that appear in both EID2 and our dataset rather than in only one

193    of the two, as these species are presumably more well-known (Figure 2C). Similarly, while

194    we believe that our model was successful at accounting for variation in host-level diversity

195    and study effort that influences network topology (see above; Figure SI1), there are certain

196    inherent biases in the training data which must be considered when interpreting our findings.

197    Most notably, viral sharing estimates in our dataset may be affected by the fact that zoonotic

198    discovery efforts commonly search limited geographic regions for a specific virus or group of

199    viruses, artificially increasing the likelihood of detecting these viruses in the same region

200    compared to a geographically random sampling regime. Moreover, when a mammal species

201    (e.g., a bat) is found with a focal virus (e.g., an ebolavirus), it is logical for researchers to then

202    investigate similar, closely related species in nearby locales[33]. These sampling approaches

203    could disproportionately weight the network towards finding phylogeographic effects on viral

204    sharing probability. However, it is highly encouraging that our model predicted patterns in

205    the external EID2 dataset, which was constructed using different data compilation methods

206    but also comprises global data covering several decades of research[29]. In sum, we believe that

207    our approach is a conservative method for minimising the biases inherent in the data. The

208    knowledge that the observed mammalian virome is biased ultimately calls for more uniform

209    viral sampling across the mammal class and increased coverage of rarely-sampled groups,

210    lending support to ongoing efforts to systematically catalogue mammalian viral diversity[3].

211

212

213 Figure 2: The modelled mammalian viral sharing network predicts observed viral sharing trends in an

214 independent dataset. In all figures, points are jittered along the x axis according to a density function;

215 the black points and associated error bars are means +/- standard errors. A: species pairs with higher

216 predicted viral sharing probability from our model were more likely to be observed sharing a virus in

217 the independent EID2 dataset. This comparison excludes species pairs that were also present in our

218 training data. B: species that hosted a zoonotic virus in our dataset had more viral sharing links in the

219 predicted all-mammal network than those without zoonotic viruses. C: species that had never been

220 observed with a virus have fewer links in the predicted network than species that hosted viruses in the

221 EID2 dataset only, in our training data only, or in both. The y axis represents viral sharing link

222 number, scaled to have a mean of 0 and a standard deviation of 1 within each order for clarity. Figure

223 SI4 displays these same data without the within-order scaling.

## Taxonomic and geographic patterns of predicted viral sharing

Our network predicted strong taxonomic patterns in the probability of viral sharing. Looking across mammalian orders, rodents (Rodentia) and bats (Chiroptera) had the most predicted species-level viral links, while carnivores and artiodactyl ungulates had substantially fewer (Figure 3A). Examining multiple mammalian orders allowed us to partition the predicted sharing network into within- and between-order links to investigate whether certain orders are better-connected to other orders. Indeed, this partitioning revealed differences in taxonomic and geographic patterns of viral sharing. In bats and rodents, large numbers of within-order links are driven by high within-order species diversity (Figure 3C). Interestingly, when within-order links were ignored, leaving only out-of-order links, rodents and bats were among the least-connected Eutherian orders (Figure 3E), while even-toed ungulates and carnivores were ranked among the most-connected (Figure 3E). Taken together, these results imply that while bats and rodents are important in viral sharing networks, their sharing is mainly restricted to other bats and rodents, respectively. This distinction only applied to mean link numbers; when link numbers were summed, rodents and bats remained highly connected regardless of which metric was used, as a result of their species richness (Figure SI5).
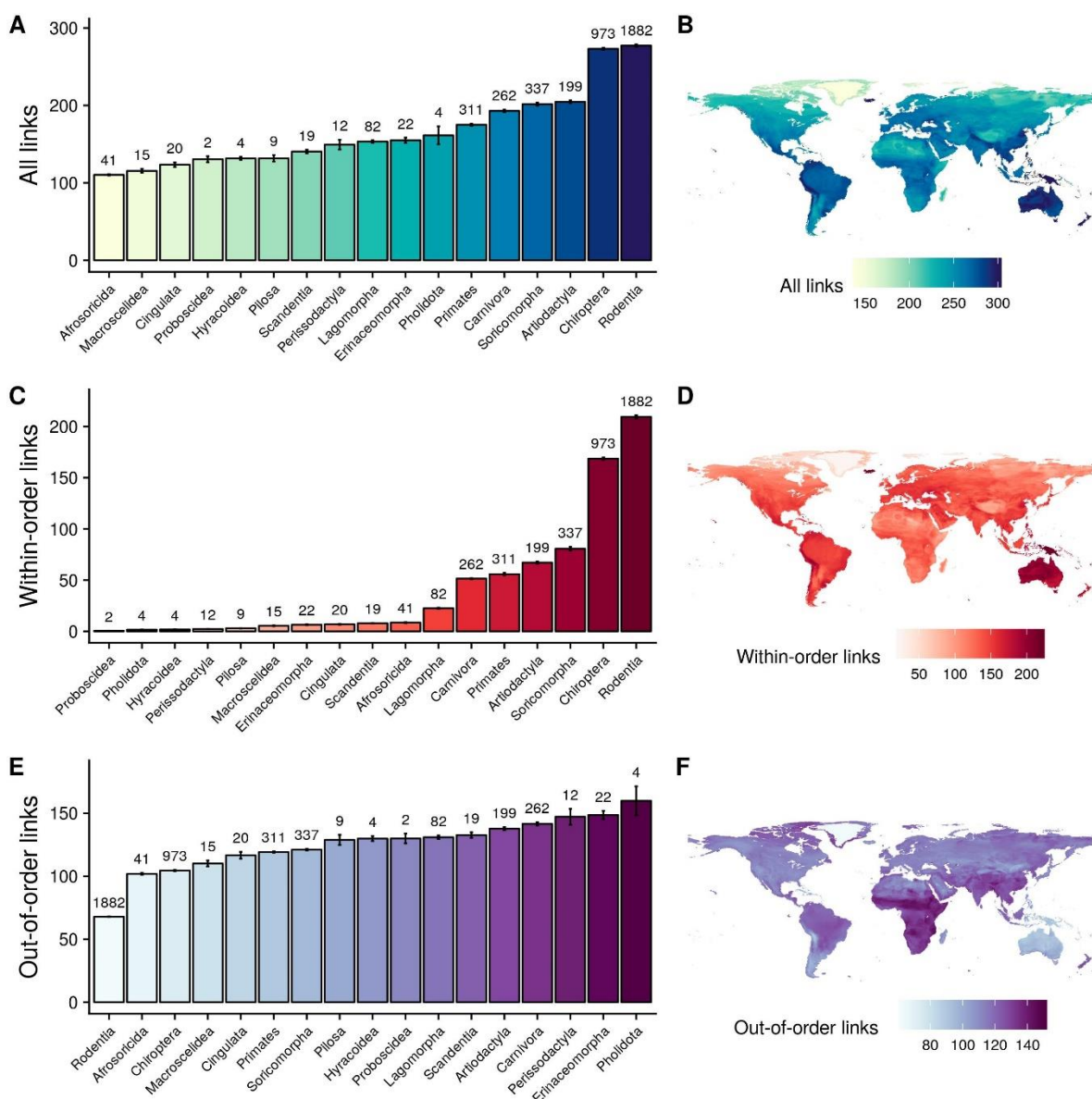
Previous analyses have demonstrated that both bats and rodents are important for hosting zoonotic viruses, with possible explanations including species-level phenotypic traits such as behaviour[5], life history[9], or metabolic idiosyncracies[34]. Our results imply that while both orders potentially host many zoonoses purely as a result of their species richness (Figure SI5), the vast majority of their viral sharing occurs within-order even though larger phylogenetic jumps are necessary for spillover. Intriguingly, recent work has shown that infection of an aggregated phylogenetic selection of hosts is an important contributor to viral zoonotic potential[35]. Rodents' and bats' tendency towards high viral interconnectedness could encourage viruses to achieve such aggregation, leading to opportunities for spillover into humans. In our analysis, both orders' high centrality emerged purely as a result of their phylogenetic diversity and geographic distributions, rather than from other phenotypic traits. If well-connected species in our network are more likely to maintain a high diversity of viruses (e.g., via multi-host dynamics leading to an expanded threshold population size[36]), this may contribute to the high viral diversity documented in bats and rodents[5]. Efforts to prioritise viral sampling regimes should consider biogeography and mammal-mammal

257    interactions in addition to searching for species-level traits associated with high viral

258    diversity.

259

260    Encouragingly, our network showed predictable scaling laws similar to those of other known

261    ecological networks[37]. Viral link numbers in within-order subnetworks (e.g., between

262    different bat species) correlated strongly with species diversity within each order ($R^2$=~0.85),

263    following a power law with a Z value of ~0.8 (Figure SI6). Similarly, out-of-order links (e.g.,

264    between a bat and a rodent) scaled linearly with the product of the species richness of both

265    orders (Figure SI7).

266



267

268    Figure 3: Taxonomic and geographic patterns of mean predicted viral sharing link numbers (degree centrality).

269      Top row: all viral sharing links; middle row: viral sharing links with species in the same order; bottom row:

270     viral sharing links with species in another order. A,C,E: average species-level viral sharing link numbers for

271     mammalian orders in our dataset. Bars represent means; error bars represent standard errors. B,D,F: geographic

272     distributions of mean viral sharing link numbers. Distributions were derived by summing the viral sharing link

273     numbers of all species inhabiting a 25km$^2$ grid square and dividing them by the number of species inhabiting the

274     grid square, giving mean degree number at the grid level.

275

276    To visualize geographic patterns of viral sharing, we projected species-level degree centrality

277    across the species' ranges then calculated grid cell-level mean degree centrality (Figure 3B),

278    as well as summed degree centrality (Figure SI5). Average centrality peaked in tropical areas

279    of South and Central America, Sub-Saharan Africa, and Southeast Asia, especially in the

280    Andes and Himalayas (Figure 3B). These patterns align with previously-reported hotspots of

281    emerging zoonoses and predicted viral diversity[5,32] and imply that areas of high biodiversity

282    are centres of viral sharing not just because of the number of overlapping species (i.e., high

283    species richness), but also because more closely related species create a more connected viral

284    sharing network in these areas. This densely-connected network structure and the increased

285    biomass present in the tropics might have synergistic implications for cross-species

286    maintenance and transmission of viral diversity in these areas. The geographic distributions

287    of mean predicted within- and between-order viral links differed notably from the distribution

288    of interspecific links generally: the relative importance of South America and East Asia was

289    higher for within-order links (Figure 3D), while Sub-Saharan Africa remained a hotspot for

290    out-of-order links (Figure 3F). Geographic patterns of summed link numbers more closely

291    mirrored underlying host species richness, whether for all links, within-order links, or out-of-

292    order links (Figure SI5).

293

294    We acknowledge that our phylogeographic model of viral sharing does not account for

295    complex ecological interactions such as coinfection or coevolution, which could impact how

296    patterns of exposure and host susceptibility translate to realised viral diversity. Future

297    investigations could extend our framework to simulate the dynamic co-speciation of

298    mammals and their viruses in order to account for these processes and/or to explicitly

299    investigate how viral sharing connectivity and viral diversity are correlated across mammal

300    species. Our model may also prove useful for building and parameterising much-needed

301    multi-host network models for conservation purposes, particularly where there is scarce prior

302    information on interspecific pathogen sharing[36,38].

303

## The network as a predictive tool

Identifying potential hosts for known and novel viruses is an important component of preemptive zoonotic disease surveillance that can speed public health responses. Predictive techniques based on species-level phenotypic and genomic data have been suggested to help prioritise sampling targets[6,7,9]. Although these approaches represent a promising methodological advance, they may not elucidate the mechanistic underpinnings of viral host range, reducing their potential efficacy for guiding public health interventions. In addition, genomic approaches require viral sequence data, which can be time-consuming and operationally challenging to acquire or share publicly. We therefore interrogated our predicted viral sharing network to investigate whether it could be used to identify potential hosts of known viruses at the species level. Using a leave-one-out prediction process (see Methods), our model showed a surprisingly strong ability to predict observed host species for 250 viruses with at least two known (non-human) mammal hosts.

We investigated the predictive potential of our model by iteratively selecting all but one of the known hosts for a given virus, then using the predicted sharing patterns of the remaining hosts to identify how the focal (removed) host was ranked in terms of its sharing probability. In practical terms, these species-level rankings could set sampling priorities for public health efforts seeking to identify hosts of a novel zoonotic virus, where one or more hosts are already known. Across all 250 viruses, the median ranking of the left-out host was 72 out of a potential 4196 mammals (i.e., in the top 1.7% of potential hosts). To compare this ranking to alternative heuristics, we examined how high the focal host would be ranked using simple ranked phylogenetic relatedness or spatial overlap values alone (i.e., the most closely-related, followed by the second-most-related, etc.). Using this method, the focal host was ranked 288th (for phylogeny) or 283rd (for space), identifying the focal host in the top 7% of potential hosts and demonstrating that sampling prioritization schemes based on our phylogeographic model would require only ¼ as many sampling targets in order to identify the correct sharing host. Our model therefore represents a substantial improvement over search methods that involve only spatial or phylogenetic similarity. Our model performed similarly at identifying focal hosts in the EID2 dataset[29]: for the 109 viruses in the EID2 dataset with more than one host, the focal host was identified in the top 63 (1.5%) potential hosts. In contrast, ranked spatial overlap predicted the focal host in the top 560 hosts, and phylogenetic relatedness in the top 174.

337

338     We observed substantial variation in our model's ability to predict known hosts among

339     different viruses. For example, the correct host was predicted first in every iteration for 7

340     viruses and in the top 10 hosts for 42 viruses. Results for 128 viruses had the focal host

341     falling within the top 100 guesses, and for only 6 viruses were the model-based host searches

342     worse than chance (focal host ranked lower than 50% of all mammals in terms of sharing

343     probability). We used this measure of viral sharing "predictability" to investigate whether

344     certain viral traits affected the ease with which phylogeography predicted their hosts. Viruses

345     with broad host phylogenetic ranges, most notably Ebola virus, challenge reservoir prediction

346     efforts since many more species must often be sampled before identifying the correct host(s).

347     To investigate whether the predictive strength of our model was limited for viruses with

348     broad host ranges and/or other viral traits, we fitted a linear mixed model (LMM) which

349     showed a strong negative association between viruses' known phylogenetic host breadth and

350     the predictability of focal hosts (model $R^2$=0.70; host breadth $R^2$=0.67; Figure SI9). This

351     association demonstrates, unsurprisingly, that predicting the hosts of generalist viruses is

352     intrinsically difficult using our method. This adds a potential limitation to the applicability of

353     our network approach, given that zoonotic viruses commonly exhibit wide host ranges[2,5]. A

354     family-level random effect accounted for little of the apparent variance in predictability

355     among viral families (Figure SI8).

356

357     Once viral host range was accounted for, hosts of vector-borne viruses were slightly easier to

358     predict than non-vector-borne viruses ($R^2$=0.1; Figure SI9) – perhaps because the sharing of

359     vector-borne viruses depends more heavily on host geographic distributions (Figure SI3).

360     Despite additional variation in the data, no other viral traits (e.g., RNA vs. DNA, segmented

361     vs. non-segmented) were important in the LMM. This implies that host phylogeographic

362     traits are a good broad-scale indicator of viral sharing, particularly when ecological specifics

363     of the virus itself are unknown.

364

## Conclusion

366     In summary, we present a simple, highly interpretable model that predicted a substantial

367     proportion of viral sharing across mammals and is capable of identifying species-level

368     sampling priorities for viral surveillance and discovery. It is worth noting that the analytical

369  framework and validation we describe were conducted on a global scale, while many
370  zoonotic sampling efforts occur on a national or regional scale. Restricting the focal
371  mammals to a regional pool may improve the applicability of our model in certain sampling
372  contexts, and future studies could leverage higher-resolution phylogenetic and geographic
373  data to fine-tune predictions. In particular, the mammalian supertree[28] has relatively poor
374  resolution at the species tips such that relatedness estimates based on alternative molecular
375  evidence (e.g., full host genome data) may allow more precise estimates of the phylogenetic
376  relatedness effect on viral sharing. Alternatively, our model could be augmented with
377  additional host, virus, and pairwise traits, using similar pairwise formulations of viral sharing
378  as a response variable, to identify ecological specificities that are critical for the transmission
379  of certain viruses, to partition viral subtypes, and, ultimately, to increase the accuracy of host
380  prediction. By generalising the spatial and phylogenetic processes that drive viral sharing, our
381  model serves as a useful guide for the prioritization of viral sampling, presenting a baseline
382  for future modelling efforts to compare against and improve upon.

383

384  Our ability to model and predict macroecological patterns of viral sharing is important in an
385  era of rapid global change. Under all conceivable global change scenarios, many mammals
386  will shift their geographic ranges, whether of their own volition or through human assistance.
387  Mammalian parasite communities will likely undergo considerable rearrangement as a result,
388  with potentially far-reaching ecological consequences[39–41]. Our findings suggest that novel
389  species encounters will provide opportunities for interspecific viral transmission, which could
390  be facilitated by even relatively small changes in range overlap. These future cross-species
391  transmission events will have profound implications for conservation and public health,
392  potentially devastating populations of host species without evolved resistance to the novel
393  viruses (e.g., red squirrel declines brought about by parapoxvirus infections spread by
394  introduced grey squirrels[42]) or increasing zoonotic disease risk by introducing viruses to
395  human-adjacent amplifier hosts (e.g., horses increasing the risk of human infection with
396  Hendra virus[20]). Thus, our global model of mammalian viral sharing provides a crucial
397  complement to ongoing work modelling the spread of hosts, vectors, and their associated
398  diseases as the result of climate change-induced range expansions[25,39].

# Acknowledgements
399

# Author contributions

411   GFA, EAE, NR, and KJO designed the study together. GFA conducted the analyses under

412   NR's supervision. GFA wrote the manuscript, while EAE, NR, and KJO offered comments

413   and edits to the manuscript throughout.

# References

415   1.   Woolhouse, M. E. J. & Gowtage-Sequeria, S. Host range and emerging and reemerging
416        pathogens. *Emerg. Infect. Dis.* **11**, 1842–7 (2005).

417   2.   Johnson, C. K. *et al.* Spillover and pandemic properties of zoonotic viruses with high host
418        plasticity. *Sci. Rep.* **5**, 1–8 (2015).

419   3.   Carroll, D. *et al.* The Global Virome Project. *Science (80-. ).* **359**, 872–874 (2018).

420   4.   Carlson, C. J., Zipfel, C. M., Garnier, R. & Bansal, S. Global estimates of mammalian viral
421        diversity accounting for host sharing. *Nat. Ecol. Evol.* **3**, 1070–1075 (2019).

422   5.   Olival, K. J. *et al.* Host and viral traits predict zoonotic spillover from mammals. *Nature* **546**,
423        646–650 (2017).

424   6.   Han, B. A. *et al.* Undiscovered Bat Hosts of Filoviruses. *PLoS Negl. Trop. Dis.* **10**, e0004815
425        (2016).

426   7.   Babayan, S. A., Orton, R. J. & Streicker, D. G. Predicting reservoir hosts and arthropod
427        vectors from evolutionary signatures in RNA virus genomes. *Science (80-. ).* **362**, 577–580
428        (2018).

429   8.   Luis, A. D. *et al.* A comparison of bats and rodents as reservoirs of zoonotic viruses: are bats
430        special? *Proc. R. Soc. B Biol. Sci.* **280**, 20122753 (2013).

431   9.   Han, B. A., Schmidt, J. P., Bowden, S. E. & Drake, J. M. Rodent reservoirs of future zoonotic
432        diseases. *Proc. Natl. Acad. Sci.* **112**, 7039–7044 (2015).

433   10.  Plourde, B. T. *et al.* Are disease reservoirs special? Taxonomic and life history characteristics.

434   *PLoS One* **12**, e0180716 (2017).

435 11. Dallas, T. A. *et al.* Host traits associated with species roles in parasite sharing networks. *Oikos*
436   **128**, 23–32 (2019).

437 12. Plowright, R. K. *et al.* Pathways to zoonotic spillover. *Nat. Rev. Microbiol.* **15**, 502–510
438   (2017).

439 13. Ge, X. Y. *et al.* Isolation and characterization of a bat SARS-like coronavirus that uses the
440   ACE2 receptor. *Nature* **503**, 535–538 (2013).

441 14. Huang, S., Bininda-Emonds, O. R. P., Stephens, P. R., Gittleman, J. L. & Altizer, S.
442   Phylogenetically related and ecologically similar carnivores harbour similar parasite
443   assemblages. *J. Anim. Ecol.* **83**, 671–680 (2014).

444 15. Wells, K. *et al.* Global spread of helminth parasites at the human-domestic animal-wildlife
445   interface. *Glob. Chang. Biol.* **24**, 3254–3265 (2018).

446 16. Stephens, P. R. *et al.* Parasite sharing in wild ungulates and their predators: Effects of
447   phylogeny, range overlap, and trophic links. *J. Anim. Ecol.* **88**, 1017–1028 (2019).

448 17. Streicker, D. G. *et al.* Host Phylogeny Constrains Cross-Species Emergence and Establishment
449   of Rabies Virus in Bats. *Science (80-. ).* **329**, 676–679 (2010).

450 18. Willoughby, A. R., Phelps, K. L., Predict Consortium,  predict@ucdavis edu & Olival, K. J. A
451   comparative analysis of viral richness and viral sharing in cave-roosting bats. *Diversity* **9**, 1–16
452   (2017).

453 19. Davies, T. J. & Pedersen, A. B. Phylogeny and geography predict pathogen community
454   similarity in wild primates and humans. *Proc. R. Soc. B Biol. Sci.* **275**, 1695–1701 (2008).

455 20. Glennon, E. E. *et al.* Domesticated animals as hosts of henipaviruses and filoviruses: A
456   systematic review. *Vet. J.* **233**, 25–34 (2018).

457 21. Chua, K. B. *et al.* Nipah Virus: A Recently Emergent Deadly Paramyxovirus. *Science (80-. ).*
458   **288**, 1432–1435 (2000).

459 22. Guy, C., Thiagavel, J., Mideo, N. & Ratcliffe, J. M. Phylogeny matters: Revisiting 'a
460   comparison of bats and rodents as reservoirs of zoonotic viruses'. *R. Soc. Open Sci.* **6**, 181182
461   (2019).

462 23. Pybus, O. G., Tatem, A. J. & Lemey, P. Virus evolution and transmission in an ever more
463   connected world. *Proc. R. Soc. B Biol. Sci.* **282**, 20142878 (2015).

464 24. Sanjuán, R., Nebot, M. R., Chirico, N., Mansky, L. M. & Belshaw, R. Viral mutation rates. *J.*
465   *Virol.* **84**, 9733–9748 (2010).

466 25. Ryan, S. J., Carlson, C. J., Mordecai, E. A. & Johnson, L. R. Global expansion and
467   redistribution of Aedes-borne virus transmission risk with climate change. *PLoS Negl. Trop.*
468   *Dis.* **13**, e0007213 (2019).

469 26. Drake, J. M. & Beier, J. C. Ecological niche and potential distribution of *Anopheles arabiensis*
470   in Africa in 2050. *Malar. J.* **13**, 213 (2014).

471 27. IUCN. The IUCN Red List of Threatened Species. *IUCN Red List of Threatened Species.*
472   *Version 2019-2* (2019). Available at: https://www.iucnredlist.org.

473 28. Fritz, S. A., Bininda-Emonds, O. R. P. & Purvis, A. Geographical variation in predictors of
474   mammalian extinction risk: big is bad, but only in the tropics. *Ecol. Lett.* **12**, 538–549 (2009).

475 29. Wardeh, M., Risley, C., Mcintyre, M. K., Setzkorn, C. & Baylis, M. Database of host-
476   pathogen and related species interactions, and their global distribution. *Sci. Data* **2**, 150049

477        (2015).

478   30.   Gómez, J. M., Nunn, C. L. & Verdú, M. Centrality in primate-parasite networks reveals the
479         potential for the transmission of emerging infectious diseases to humans. *Proc. Natl. Acad. Sci.*
480         **110**, 7738–41 (2013).

481   31.   Wang, Y. X. G. G. *et al.* Phylogenetic structure of wildlife assemblages shapes patterns of
482         infectious livestock diseases in Africa. *Funct. Ecol.* **33**, 1332–1341 (2019).

483   32.   Allen, T. *et al.* Global hotspots and correlates of emerging zoonotic diseases. *Nat. Commun.* **8**,
484         1124 (2017).

485   33.   Becker, D. J., Crowley, D. E., Washburne, A. D. & Plowright, R. K. Temporal and spatial
486         limitations in global surveillance for bat filoviruses and henipaviruses. *Biol. Lett.* **15**,
487         20190423 (2019).

488   34.   Xie, J. *et al.* Dampened STING-Dependent Interferon Activation in Bats. *Cell Host Microbe*
489         **23**, 297–301 (2018).

490   35.   Park, A. W. Phylogenetic aggregation increases zoonotic potential of mammalian viruses. *Biol.*
491         *Lett.* **15**, 20190668 (2019).

492   36.   Lloyd-smith, J. O. *et al.* Epidemic Dynamics at the Human-Animal Interface. *Science (80-. ).*
493         **326**, 1362–1368 (2009).

494   37.   Brose, U., Ostling, A., Harrison, K. & Martinez, N. D. Unified spatial scaling of species and
495         their trophic interactions. *Nature* **108**, 167–171 (2003).

496   38.   Silk, M. *et al.* Integrating social behaviour, demography and disease dynamics in network
497         models: applications to disease management in declining wildlife populations. *Philos. Trans.*
498         *R. Soc. B* **374**, 20180211 (2019).

499   39.   Carlson, C. J. *et al.* Parasite biodiversity faces extinction and redistribution in a changing
500         climate. *Sci. Adv.* **3**, e1602422 (2017).

501   40.   Williams, J. E. & Blois, J. L. Range shifts in response to past and future climate change: Can
502         climate velocities and species' dispersal capabilities explain variation in mammalian range
503         shifts? *J. Biogeogr.* **45**, 2175–2189 (2018).

504   41.   Chen, I.-C., Hill, J. K., Ohlemüller, R., Roy, D. B. & Thomas, C. D. Rapid range shifts of
505         species associated with high levels of climate warming. *Science (80-. ).* **333**, 1024–6 (2011).

506   42.   Tompkins, D. M., Sainsbury, A. W., Nettleton, P., Buxton, D. & Gurnell, J. Parapoxvirus
507         causes a deleterious disease in red squirrels associated with UK population declines. *Proc. R.*
508         *Soc. B Biol. Sci.* **269**, 529–533 (2002).

509   43.   R Core Team. R: A language and environment for statistical computing. R Foundation for
510         Statistical Computing, Vienna, Austria. (2018).

511   44.   Wood, S. N. Fast stable restricted maximum likelihood and marginal likelihood estimation of
512         semiparametric generalized linear models. *J. R. Stat. Soc. Ser. B (Statistical Methodol.* **73**, 3–
513         36 (2011).

514

# Methods

515

**Making the training data network**

516

517    Code and data for all analyses are available at

518    https://github.com/gfalbery/ViralSharingPhylogeography. Our dataset included 1920 mammal-

519    virus associations obtained from an exhaustive literature search which has been used to

520    investigate how species traits influence mammalian viral diversity[5]. We removed humans and

521    rabies virus from the dataset as both were disproportionately well-connected, and we

522    removed 20 non-Eutherian mammals because they were extreme phylogenetic outliers,

523    leaving 591 Eutherian mammals that shared 401 viruses. We made an unweighted bipartite

524    network using the mammal-virus associations and projected the unipartite mammal-mammal

525    network, which we then converted into a sequence of all unique mammal-mammal pairs

526    where 1/0 denoted whether the pair of species shared a virus or not. This comprised only the

527    lower triangle of the adjacency matrix to avoid duplicating associations and to remove self-

528    connections, and only included mammals with at least one sharing link (final N=174345

529    unique mammal-mammal pairs). 6.4% of these pairs shared at least one virus.

530

531    All analyses were performed in R version 3.6.0[43]. Phylogenetic similarity was calculated

532    using a mammalian supertree[28] as previously described[5]. Pairwise phylogenetic distances

533    were defined as the cumulative branch length between the two species and were scaled to

534    between 0 and 1, and subtracted from 1 to give a measure of relative phylogenetic similarity

535    (rather than distance). Of the 4716 Eutherian species in the mammalian supertree, 591 had

536    virus association records in our fully-connected network and 4196 had known geographic

537    ranges. We used IUCN species ranges to quantify species' geographic distributions[27]. These

538    range maps are generated based on expert knowledge and only comprise species

539    presence/absence information rather than density. We converted all range polygons to $25 \text{ km}^2$

540    raster grids. For each species-pair, we quantified range overlap as the number of raster grid

541    squares jointly inhabited by the two species (in the Mollweide projection, which exhibits

542    equal grid size), divided by the total number of grid squares occupied by these species

543    combined, so that each value was scaled from 0-1: $\text{overlap}_{A,B} = \text{grid}_{A,B}/(\text{grid}_A + \text{grid}_B - \text{grid}_{A,B})$.

544    Disease-related research effort for each host species was quantified as previously described,

545    using counts of studies including species names and disease-related terms such as "virus,"

546    "pathogen", or "parasite"[5]. To fit citation number as a pairwise trait, we took the smaller of a

547    pair of species' respective citations, and log-transformed the value. Domestication status was

548    defined *sensu lato*, again as previously described[5], based on whether a species was ever seen

549    in a domestic setting. We fit this as a binary pairwise trait where 1=at least one of the species

550    was domesticated and 0=neither species had been domesticated.

551 **Model formulation**

552 We fitted a Generalised Additive Mixed Model (GAMM) to examine which traits influenced

553 viral sharing among mammal pairs using accelerated discretized implementation in the **mgcv**

554 package[44]. We fitted viral sharing (0/1) as the response variable, with a binomial family

555 specification. The model had the following structure:

556

557 Bernoulli(Viral sharing) ~ s(Phylogenetic similarity, by = ordered(Gz)) +

558       t2(Phylogenetic similarity, Geographic overlap, by = ordered(!Gz)) +

559       Minimum citation number + Domestication status +

560       mm (Species 1 + Species 2)

561

562 The first term ("s") represents a phylogeny effect smooth fitted across species pairs that did

563 not overlap in space (Gz=1), and "t2" represents a phylogeny:geography tensor product

564 smooth fitted to species that had geographic overlap greater than zero (Gz=0). This allowed

565 us to model these two aspects of the data separately, helping us to more effectively model the

566 large number of spatial zeroes (85% of species pairs did not overlap in space). "mm"

567 represents a multi-membership random effect, accounting for the identity of both species in

568 the pair. We implemented this multi-membership effect to control for species-level effects by

569 including a species-level effect for both the row (Species 1) and column (Species 2) of the

570 sharing matrix. Using the paraPen specification in **mgcv**, these random effects were

571 constrained to sample from the same distribution, resulting in a single estimate of the

572 variance associated with each unique species. Most precisely, these effects in our model help

573 capture variation in viral sharing that could likely be explained by species-level factors that

574 are unobserved or otherwise excluded (i.e., differences in underlying viral diversity, which

575 would be expected to positively impact the probability of interspecific sharing). In sum, this

576 model formulation allowed us to estimate the effect of pairwise predictors (geographic

577 overlap, phylogenetic similarity) in determining viral sharing as well as evaluate the

578 influence of species identity.

579

580 To investigate whether the effects of geography and phylogeny depended on which subset of

581 viruses we investigated, we fit the model to non-exclusive subnetworks of mammal-mammal

582 pairs based on the types of viruses they were connected by. Viral subtypes included RNA

583 viruses (566 hosts sharing 381 viruses); vector-borne RNA viruses (333 hosts sharing 164

584   viruses); non-vector-borne RNA viruses (391 hosts sharing 205 viruses); and DNA viruses

585   (151 hosts sharing 205 viruses). There were only 2 vector-borne DNA viruses in our data. We

586   eliminated from each analysis any hosts that were not carrying the focal virus type.

587

588   We elected to use a binary model of viral sharing (0/1) rather than an integer count model

589   (0+) for two reasons: first, the data distribution was highly skewed, with few very large

590   values and many zeroes. Under these conditions, we found a count-based model formulation

591   including species-level random effects computationally intractable. Second, the observed

592   viral diversity is likely a considerable underestimate, but the extent of this underestimate is a

593   matter of hot debate[3,4]. As such, the predictions for viral sharing from such a model could be

594   relative and biased, while binary models offer a more appropriate resolution to quantify

595   sharing patterns. We wished to avoid estimating a precise number of viruses shared among

596   pairs of species for this reason.

597

598   **Model validation**

599   To check the fit of the model, we predicted 0/1 viral sharing values from the model 1000

600   times and examined how the values compared to the proportions of 0's and 1's in the

601   observed data, finding high agreement between the two. We repeated this procedure using a)

602   the full dataset; b) only the fixed effects, with random effects randomised in each iteration;

603   and c) only the random effects, with fixed effects held at the mean. We then used these

604   predicted links to create 1000 unipartite viral sharing networks, estimating link numbers

605   (degree centrality) for species in each replicated network. We took the mean of these values

606   across the 1000 replicated networks to give the predicted values displayed in Figure SI1.

607

608   We quantified deviance contributions of our explanatory variables by calculating model

609   deviance when dropping each variable, and comparing these against the full model and an

610   intercept-only model deviance. For each of our explanatory variables (geographic overlap,

611   phylogenetic similarity, minimum citation number, domestication status, and species-level

612   random effects) we randomised the observed values 1000 times, then predicted sharing

613   probabilities for these values using our model estimates. This randomisation procedure

614   allowed us to predict while accounting for the uneven data distribution, rather than using

615   mean values.

**Simulating viral sharing networks**

Following reconstruction of the observed network as part of our model validation, we repeated the prediction process on an exhaustive mammal dataset to estimate viral sharing across all mammals. We set minimum citation number to the data mean, and set domestication status to 0. We repeated the predictions 1000 times, randomising the species-level random effects each time. The full prediction dataset included 4196 Eutherian mammals with known spatial distributions and phylogenetic associations, resulting in 8.8 million unique pairwise combinations. After predicting 1000 binary sharing networks across all mammals, we summarised the average predicted link number (degree centrality) of each species across the 1000 replicates. We then calculated the mean species-level link number within each mammalian order to examine taxonomic patterns. To project the spatial patterns of connectedness, we assigned each species range polygon the link number (degree centrality) of its host species[27] and took the mean value for each grid square, thereby correcting for species richness. We then repeated these taxonomic and geographic summaries using within-order and between-order link numbers separately. We also took the summed values, which more closely reflect underlying patterns of species richness.
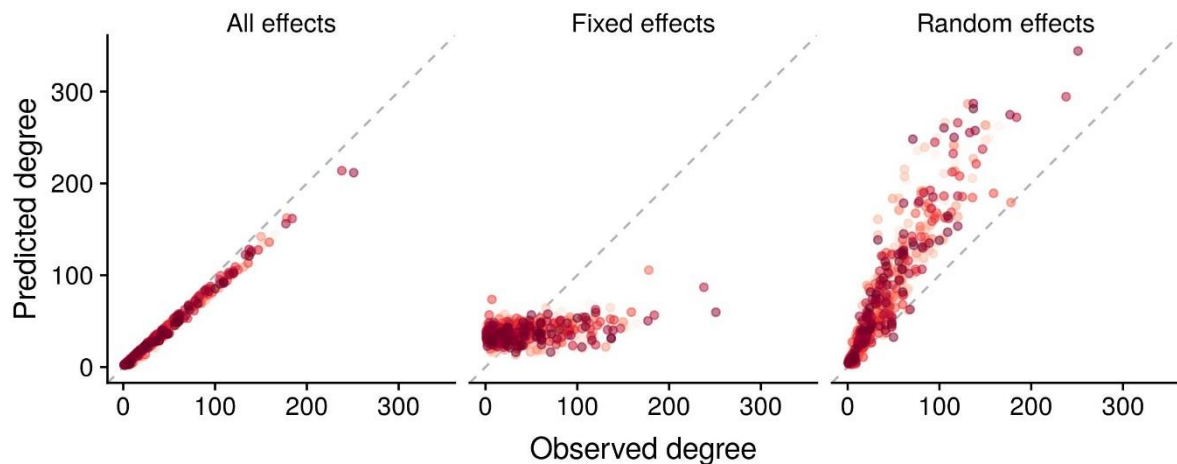
We validated the predicted network by comparing it to sharing patterns in the Enhanced Infectious Diseases Database (EID2)[29]. We eliminated species pairs that were in our training data and identified whether species pairs that shared viruses in EID2 were more likely to share viruses in our predicted network than species pairs that did not. In addition, we investigated whether species that were shown to host zoonoses in our training dataset were more highly-connected in the predicted network. Finally, we investigated whether species that were present in only EID2, in only our training data, or in both were more highly-connected in our predicted network than species that did not appear in either dataset and were therefore taken to have not been observed hosting a virus.

**Predicting hosts of focal viruses**

To investigate the ability of the model to predict known hosts of viruses in our dataset, we iteratively investigated the sharing patterns of known hosts independently for all viruses with >1 host. For each virus, we removed one host at a time, and then investigated which species the remaining known host species were likely to share viruses with based on the all-mammal predicted network. If the removed host ("focal host") was on average highly likely to share
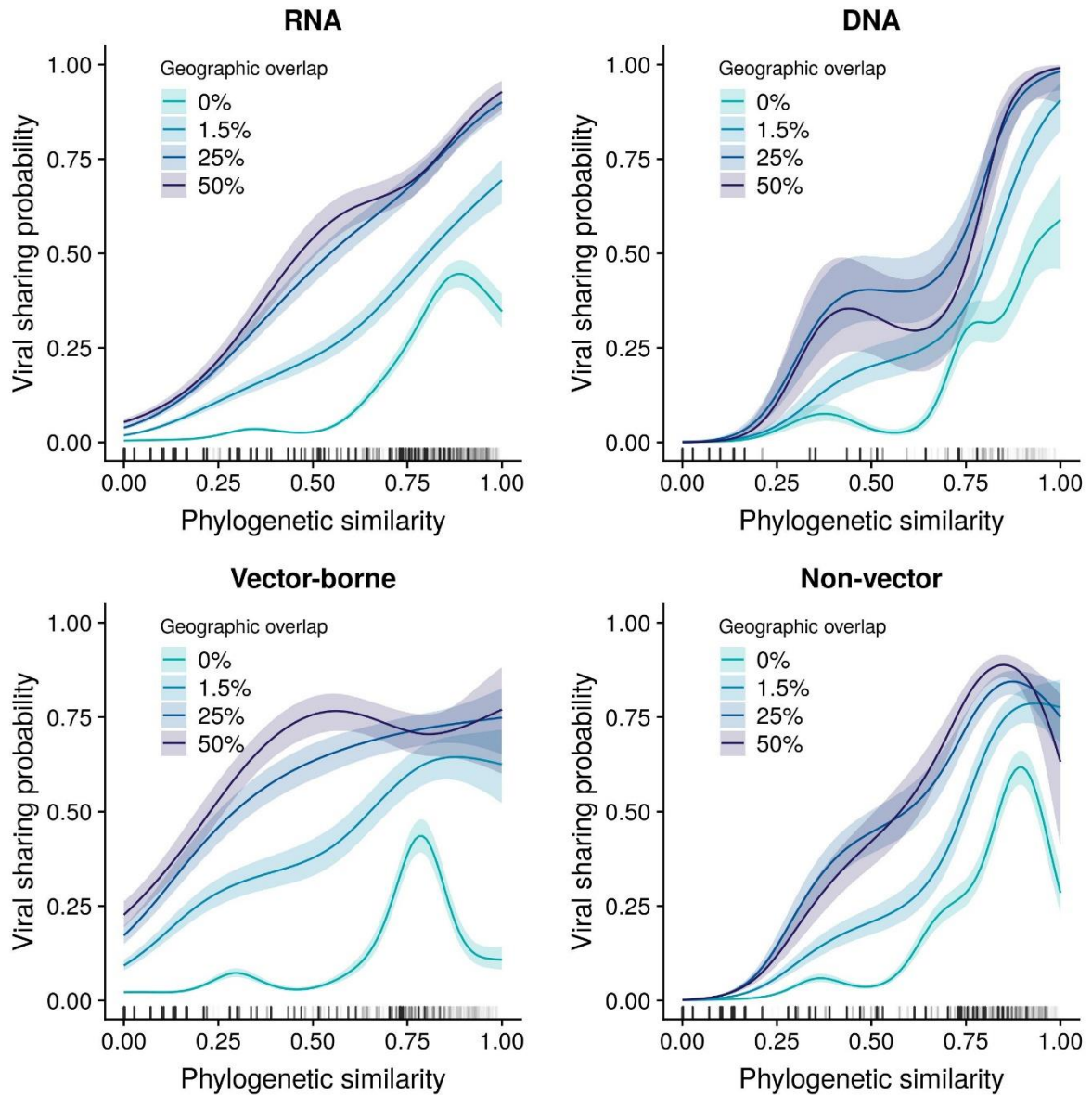
648  viruses with the remaining species, our model was taken to be useful for predicting patterns

649  of mammal sharing based on known host distributions. The mean ranking of the focal hosts

650  across each prediction iteration was used as a measure of "predictability" for each virus. We

651  carried out this process for the 250 viruses with more than one known host with associated

652  geographic and phylogenetic data and then on the 109 such viruses in the EID2 data.

653  Once the predictability of each virus was calculated, we fitted a linear mixed model

654  examining log10(mean focal host rank) as an inverse measure of predictability (higher rank

655  corresponds to decreased predictability) for each virus. We added mean phylogenetic host

656  similarity as a fixed effect and viral family as a random effect to quantify how viral

657  phylogeny affected predictability. We included additional viral traits in the model, including:

658  cytoplasmic replication (0/1); segmentation (0/1); vector-borne transmission (0/1); double- or

659  single-strandedness; DNA or RNA; enveloped or non-enveloped; or zoonotic ability (0/1 for

660  whether the virus was associated with humans in our dataset).

661



662

663  Figure SI1: Predicted degree centrality of species in our training data network, predicted using our
664  GAMM estimates. Fixed + random effects were very effective at reproducing individual species'
665  degree centrality (left); fixed effects were less effective (middle); and random effects alone had a
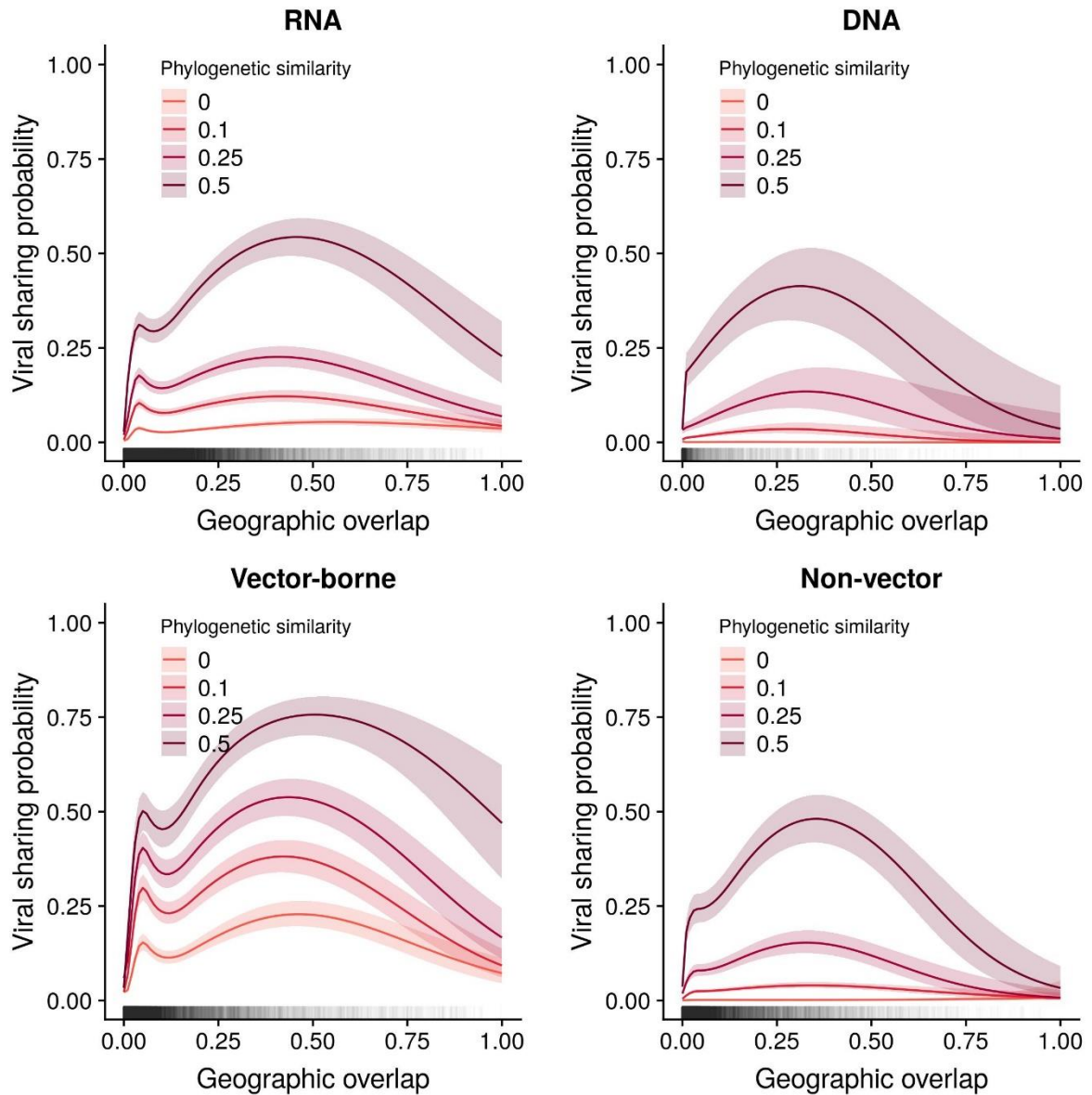666  strong but imperfect effect (right).

667

668

669 Figure SI2: GAMM-derived viral sharing estimates for the effect of phylogenetic similarity for four viral subsets
670 (top row: all RNA viruses and DNA viruses; bottom row: vector-borne RNA viruses and non-vector-borne RNA
671 viruses). Each GAMM smooth is displayed at multiple geographic overlap values (different colours).

672

673

674     Figure SI3: GAMM-derived viral sharing estimates for the effect of geographic overlap for four viral subsets
675 (top row: all RNA viruses and DNA viruses; bottom row: vector-borne RNA viruses and non-vector-borne RNA
676     viruses). Each GAMM smooth is displayed at multiple geographic overlap values (different colours).

677

678

679

| RESPONSE | SAMPLES | DEVIANCE CONTRIBUTIONS | | | | | |
|---|---|---|---|---|---|---|---|
| | | Geography | Gz | Phylogeny | Citations | Domestic | Spp |
| **ALL VIRUSES** | 591 | 0.077 | 0.067 | 0.336 | 0.005 | 0.002 | 0.512 |
| **RNA** | 566 | 0.079 | 0.067 | 0.33 | 0.005 | 0.003 | 0.516 |
| **DNA** | 151 | 0.008 | 0.031 | 0.729 | 0.001 | 0.004 | 0.227 |
| **VECTOR-BORNE** | 333 | 0.153 | 0.11 | 0.145 | 0 | 0.008 | 0.584 |
| **NON-VECTOR** | 391 | 0.011 | 0.019 | 0.625 | 0.016 | 0.001 | 0.328 |

680

681 Table SI1: The deviance contributions and sample sizes (number of hosts) for each of our viral
682 sharing GAMMs. The deviance terms are, in order: proportional geographic overlap; binary
683 geographic overlap greater than zero (0/1); phylogenetic similarity; minimum citation number;
684 domestication status; and the species-level random effect.
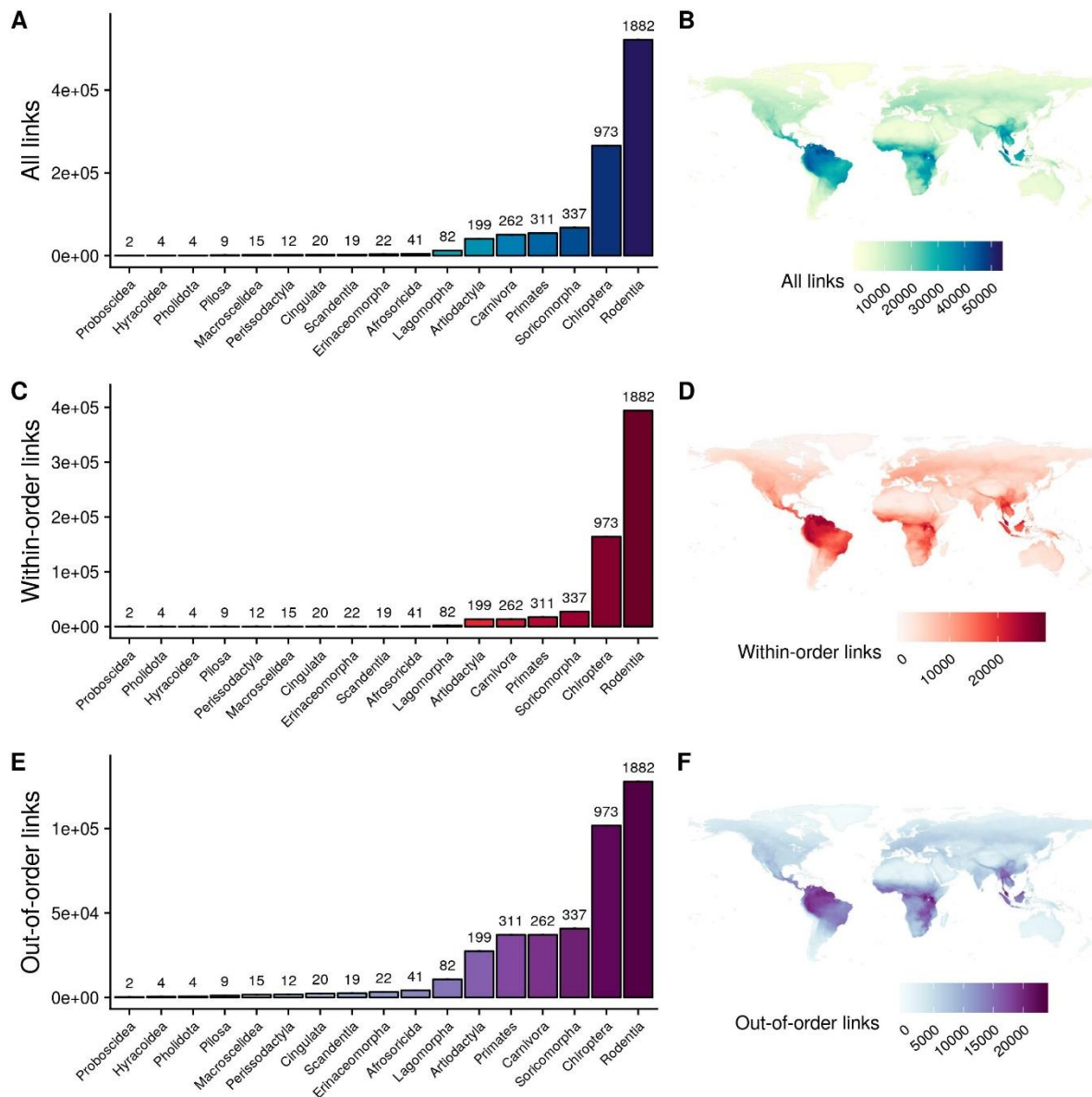
685

686

687

688



689

690 Figure SI4: Mammal species that were observed with at least one virus in the training dataset or the
691 EID2 dataset had higher degree centrality (link number) in our predicted network. This figure displays
692 the raw data that are displayed in Figure 2C in the main text, but without being scaled within orders.

693

Figure SI5: Taxonomic and geographic patterns of predicted viral link numbers. Top row: all links; middle row: links with species in the same order; bo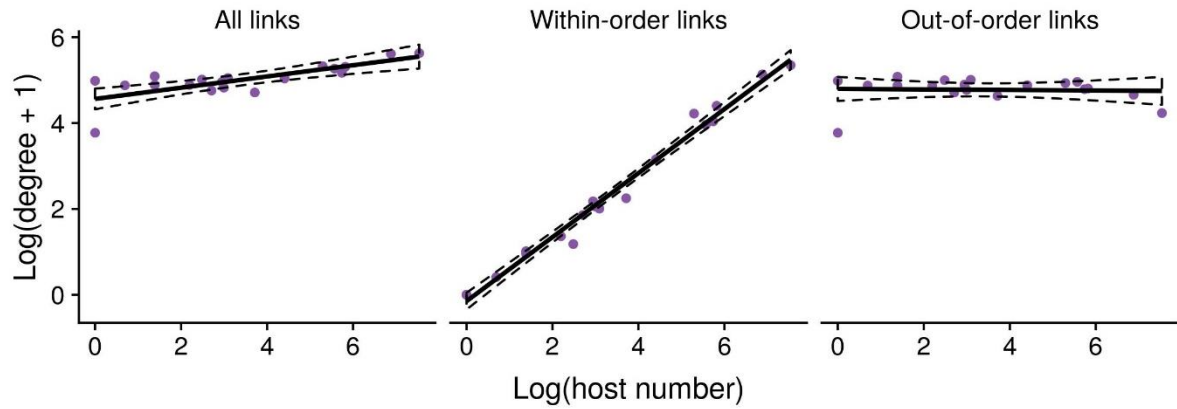ttom row: links with species in another order. A,C,E: summed species-level link numbers for mammalian orders in our dataset. B,D,F: geographic distributions of link numbers. Distributions were derived by summing the link numbers of all species inhabiting a grid square.
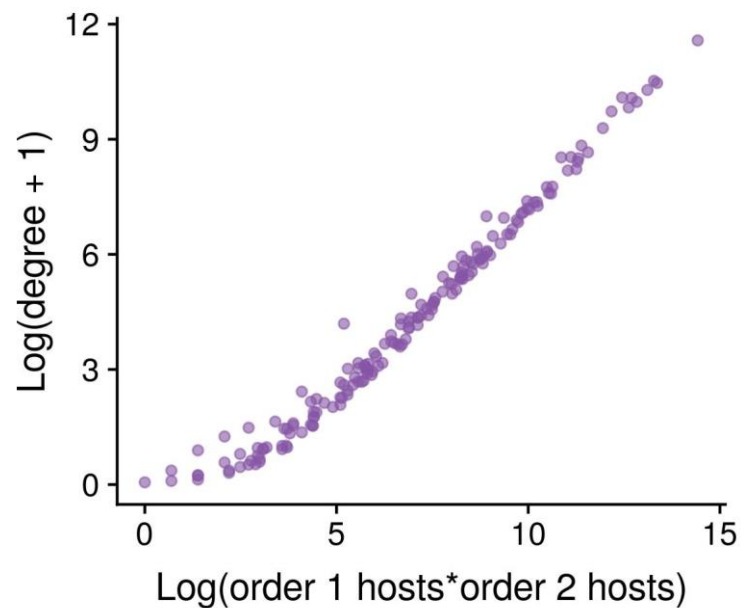
700

Figure SI6: Scaling of degree centrality (link numbers) followed a power law when looking within-orders, but not between orders. The trend line and 95% confidence intervals are derived from a linear model fitted to the data.
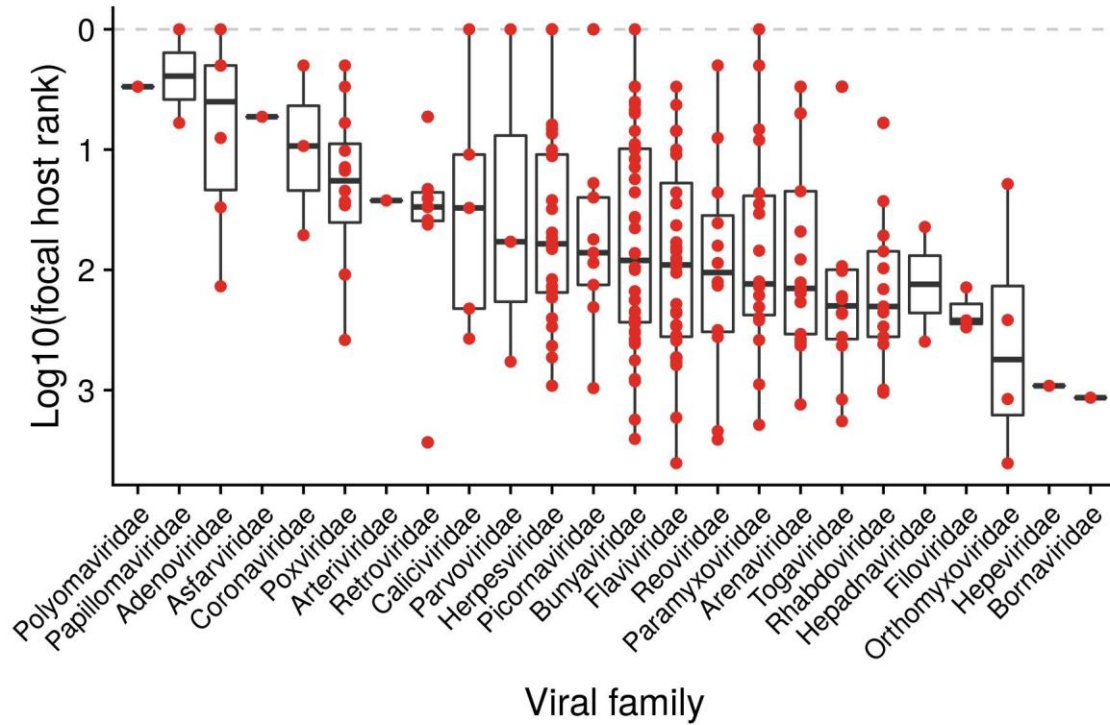
704

705



706

Figure SI7: Predicted between-order link numbers scales according to the log-product of the number of species in the two orders. Each point represents a pair of orders (N=171).
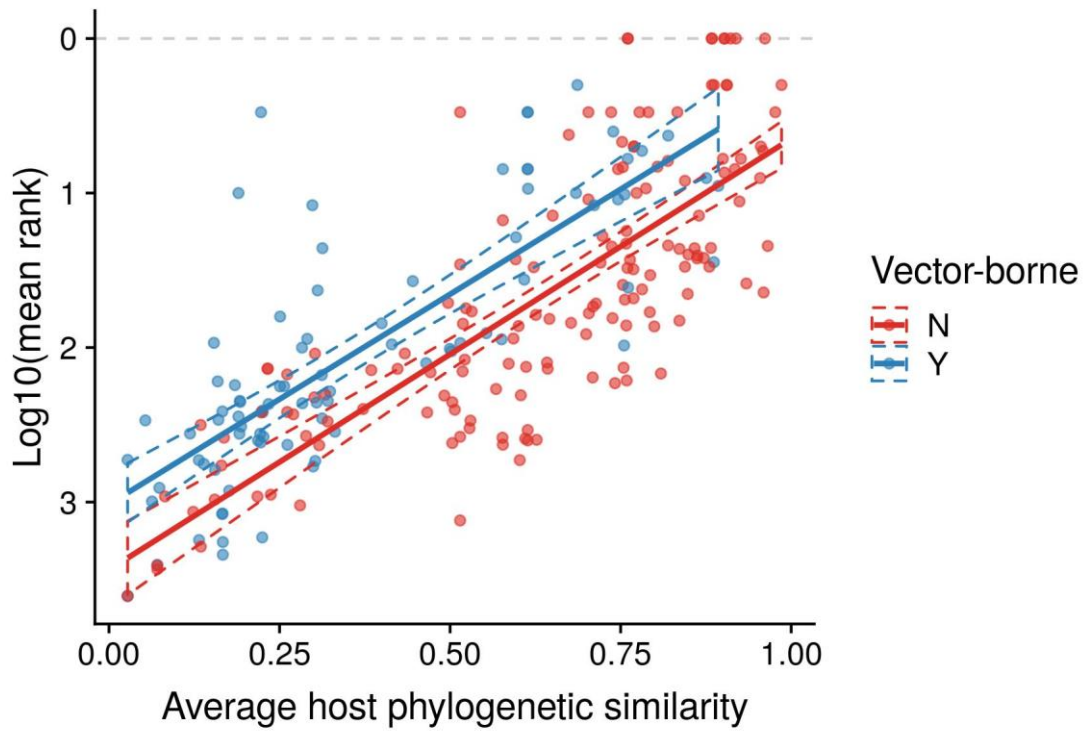
709

710

711  Figure SI8: The phylogeographic predictability of viruses' reservoir hosts varied considerably across
712  viral families, although the family-level random effect did not account for much of the model's
713  variance. Families are ordered along the x axis in order of decreasing predictability. The y axis
714  displays the mean rank of the focal host in our reservoir host prediction simulation, on a reversed
715  log10-scale. Values closer to the top of the figure represent more predictable viruses.

716

717



718

719    Figure SI9: Viral host range strongly impacted the predictability of reservoir hosts. The x axis
720    displays the mean phylogenetic similarity of a virus's hosts (i.e., an inverse measurement of viral host
721    range). The y axis displays the mean rank of the focal host in our reservoir host prediction simulation,
722    on a reversed log10-scale. Values closer to the top of the figure represent more predictable viruses.
723    The trend lines and 95% confidence intervals were derived from a linear mixed model fitted to the
724    data.

725