

Predicting the global mammalian viral sharing network using phylogeography

Gregory F Albery^{1,2,3*}, Evan A Eskew¹, Noam Ross¹, 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

Submitted to *Nature Ecology and Evolution* on 11th August 2019.

Abstract

Understanding the factors driving viral sharing among mammal species is an important public health research priority. While previous analyses have succeeded in identifying mammal taxa and phenotypic traits associated with high viral diversity, few general models have been developed to predict interspecific sharing patterns themselves. Here we show that host phylogenetic similarity and geographic range overlap are strong, non-linear predictors of viral sharing among mammal species, with effects similar in magnitude to species-level phenotypic traits. Using these traits, we predict global viral sharing patterns across 4196 mammal species and show that our simulated network successfully predicts sharing patterns and reservoir hosts using an external dataset. High rates of interspecific viral connectedness occurred in the tropics, particularly involving rodents and bats. Our results illustrate the importance of macroecological factors in shaping mammalian viral communities, providing a robust, general model to predict viral host range and guide pathogen surveillance and conservation efforts.

Introduction

Most emerging human viruses originate in wild mammals, so understanding the drivers of interspecific viral transmission in these taxa is an important public health research priority^{1,2}. Despite a rapidly expanding knowledge base, the number of identified mammalian viruses remains taxonomically biased and limited in scope, likely comprising less than 1% of the mammalian virome^{3,4}. Furthermore, host range is inadequately characterized even for the best-studied viruses⁵⁻⁷. To help prioritise pathogen discovery efforts and zoonotic disease surveillance in wildlife, studies have linked high (zoonotic) parasite diversity with certain host taxa, such as rodents and bats^{5,8}, and/or with phenotypic host traits such as reproductive output^{9,10}. Viral diversity has also been associated with host macroecological traits including geographic range size¹¹ and sympatry with other mammals⁵. The rationale for investigating viral diversity is that species with more viruses will generate more opportunities for pathogen transmission to other species, including humans. However, in order to infect a new species, a virus must transmit, invade, and potentially replicate within the novel host¹². Each of these processes becomes less likely if the two hosts differ more in terms of their geographic range, behaviour, and/or biochemistry (i.e., cellular receptors allowing viral attachment and invasion)^{12,13}. Consequently, the probability that a pair of hosts will share a virus is defined both by the species' underlying viral diversity, and by pairwise similarity measures such as spatial overlap, phylogenetic relatedness, and ecological similarity¹⁴⁻¹⁶. Previous investigations into pairwise determinants of viral sharing have been limited to single host orders (e.g., bats^{17,18}, primates¹⁹, ungulates¹⁶, and carnivores^{14,16}), while sometimes lumping together different types of parasites (e.g., helminths, viruses, and bacteria). Yet, many viruses are shared across large host phylogenetic distances (e.g., Nipah virus in bats and pigs, among many others^{20,21}), requiring a broader understanding of viral sharing across mammals to predict global patterns.

A wildlife species' ability to transmit viruses to humans (its "zoonotic potential") can be predicted by its spatial proximity and phylogenetic relatedness to humans^{5,22}, and by its centrality in mammalian parasite sharing networks²³. These observations imply that generalised phylogeographic rules for mammalian viral sharing may be applicable to humans (i.e., humans are "just another host"), opening up their potential for use in forecasting zoonotic events. Nevertheless, most research has been directed towards species' phenotypic traits and to direct links between mammals and humans (zoonoses) rather than to links among

mammal species themselves^{5,8,9}. As a result, macroecological determinants of viral sharing are less understood, and their importance relative to host-level phenotypic traits is uncertain. Here, we analyse pairwise viral sharing using a novel, conservative modelling approach designed to partition the contribution of species-level traits from pairwise phylogeographic traits. This method of analysis stands in contrast to previous studies of mammalian viral sharing which have mainly focussed on host-level traits, and importantly buffers against certain inherent biases in the observed viral sharing network, including host-species-level sampling bias, when making predictions.

Results and Discussion

Viral sharing GAMM

We fitted a Generalised Additive Mixed Model (GAMM) designed to partition the contribution of species-level effects and pairwise similarity measures to mammalian viral sharing probability. We used a published database of 1920 mammal-virus associations (excluding humans) as a training dataset⁵. These data included 591 wild mammal species, equalling 174345 pairwise host species combinations, 6.4% of which shared at least one virus. Overall, our model accounted for 44.8% of the total deviance in pairwise viral sharing, with 51.1% of this explained deviance attributed to the identities of the species involved. We included this species-level effect in our model as a multi-membership random effect capturing variation in each species' connectedness and underlying viral diversity (see methods). Our model structure was extremely effective at controlling for species-level variation in our dataset: i.e., the term had a strong impact on the centrality of each species when we simulated networks using just these parameters (Figure S11). This observation suggests that much of the dyadic structure of observed viral sharing networks is determined by uneven sampling and concentration on specific species, rather than by macroecological processes.

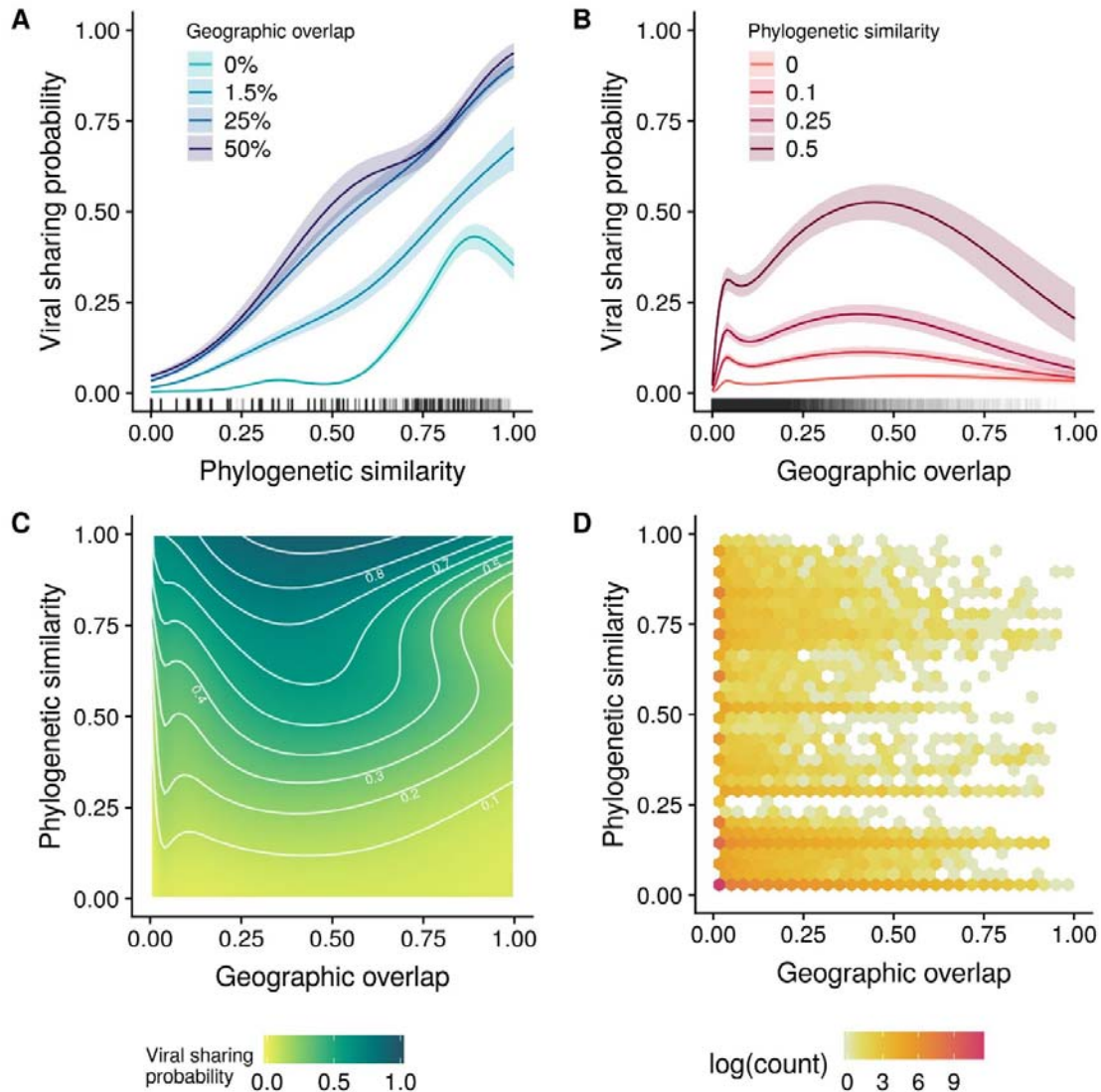


Figure 1: Viral sharing GAMM model outputs and data distribution. A: predicted viral sharing probability increases with increasing phylogenetic relatedness; the different coloured lines represent different geographic overlap values. B: predicted viral sharing probability increases with increasing geographic 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.

Increasing host phylogenetic similarity and geographic overlap were associated with increased probability of viral sharing, together accounting for the remaining 49% of explained model deviance (Figure 1A-C). Geography, phylogeny, and their interaction all showed strong non-linear effects, with geographic overlap in particular driving a rapid increase in viral sharing that began at ~0-5% range overlap values, peaked at 50% overlap values, and then levelled off (Figure 1B). Although occupying little of the visual space within the model presentation, 93% of mammal pairs had less than 5% spatial overlap (Figure 1B,D). This predicted effect therefore has implications for the vast majority of mammal-mammal interactions globally and represents an important minimum threshold of geographic overlap that allows for dramatically increased probability of interspecific viral sharing (Figure 1B). Thus, the observed viral sharing network may have been shaped substantially by initial cross-species transmission events that appear in small parts of a species range where interspecific range overlap occurs, followed by intraspecific spread. The alternative scenario is repeated interspecific transmission, which would be facilitated by larger interspecific range overlaps. Meanwhile, the slow levelling off and decreasing influence of space sharing at high values (Figure 1B) is likely driven by a few species pairs (less than 0.5% of pairs had >50% spatial overlap), so this trend may represent a statistical artefact influenced by the data distribution (Figure 1D). Importantly, the effect of geographic overlap was much stronger for species that were more closely related, but followed a qualitatively similar shape at all levels of phylogenetic relatedness (Figure 1B).

Phylogenetic similarity accounted for more than twice as much model deviance as did spatial overlap (33.8% vs 14.4%). The great majority (86%) of mammal pairs in our dataset did not overlap geographically. For these pairs, probability of viral sharing was generally close to zero and only increased when phylogenetic similarity exceeded ~0.5 (Figure 1A). This phylogenetic distance corresponds roughly to order-level similarity; that is, if two species did not overlap in space, it was highly unlikely that they shared a virus unless they were within the same taxonomic order (8% of pairs). With increasing amounts of spatial overlap, the gradient of the phylogeny effect increased and became more linear (Figure 1A). Our findings support the important role of mammalian evolutionary history in shaping contemporary patterns of viral sharing and diversity^{5,24}. The greater importance of phylogeny relative to geography contrasts with previous analyses concerning viral sharing in primates¹⁹ and ungulates¹⁶, likely reflecting the wider phylogenetic range of hosts considered here. In contrast to geography and phylogeny, minimum citation count and domestication status

accounted for a vanishingly small amount of the deviance in viral sharing probability (0.2% and 0.1%, respectively), even though they have important effects on observed viral diversity in this dataset⁵. Their impact on viral sharing was largely accounted for by the species-level random effects, as each of these effects was expected to represent individual species' viral diversity and connectedness rather than altering pairwise interactions.

We repeated our analysis on viral sharing subnetworks to assess whether the influence of spatial overlap and phylogenetic relatedness depended on viral subgroup. These groups included RNA viruses (N=566 hosts); vector-borne RNA viruses (333 hosts); non-vector-borne RNA viruses (391 hosts); and DNA viruses (151 hosts). The importance of geographic overlap varied widely across all groups of viruses (Figure SI2; Table SI1), while the influence of host phylogenetic relatedness was more consistent (Figure SI3; Table SI1). Generally, phylogeny was more important in determining sharing of DNA viruses than it was for RNA viruses, while space sharing was more important for vector-borne RNA viruses, and less so for non-vector-borne RNA viruses. These results likely reflect important aspects of viral ecology, transmission, and evolution. For example, directly-transmitted RNA viruses are fast-evolving, allowing them to more quickly adapt to novel hosts, such that phylogenetic distances are less important in determining viral sharing patterns²⁵. Conversely, DNA viruses are more evolutionarily constrained, with an evolutionary rate typically <1% that of RNA viruses, such that phylogenetic distance between hosts presents a more significant obstacle for sharing of DNA viruses²⁶. The profound importance of geographic overlap in shaping the viral sharing network for vector-borne RNA viruses (Figure SI3) likely emerges from the geographic distributions and ecological constraints placed on vectors, lending further support to efforts to model the global spread of arboviruses by predicting changes in their vectors' distributions and ecological niches^{27,28}.

Using GAMM estimates to predict sharing patterns

Previous trait-based approaches to predict viral sharing and reservoir hosts have been hindered by incomplete and inconsistent characterization of traits central to those modelling efforts. In contrast, spatial distributions and phylogenetic data are readily available and uniformly quantified for the vast majority of mammals and, as we have shown, are reliable predictors of viral sharing (>20% of deviance). Thus, we used our GAMM estimates to

predict unobserved global viral sharing patterns across 8.8 million mammal-mammal pairs using a database of geographic distributions²⁹ and a recent mammalian supertree³⁰ (see methods). The predicted network included 4196 (non-human) Eutherian mammals with available data, 591 of which were recorded with viral associations in our training data. We calculated each species' predicted degree centrality (viral link number), as the simplest and most interpretable network-derived measure of viral sharing. We identified geographic and taxonomic trends in link numbers, validated our network using an external dataset, and simulated reservoir identification to assess host predictability for focal viruses (see methods).

We confirmed that our modelled network recapitulated expected patterns of viral sharing using the Enhanced Infectious Diseases Database (EID2) as an external dataset³¹. Pairs of species that share viruses in EID2, but which were not in our dataset (see methods), had a much higher mean sharing probability in our predicted network (20% versus 5%; Figure 2A). In addition, more central species in the predicted network were more likely to have been observed with a virus, whether zoonotic (Figure 2B) or non-zoonotic (Figure 2C), implying that the predicted network accurately captured realised potential for viral sharing and zoonotic spillover. This finding concurs with similar work in primates which demonstrated that high centrality in primate-parasite networks is associated with carriage of zoonoses²³. We corroborate these findings considering all mammal-mammal viral links, not just zoonotic links, and show that for each mammalian order, species with more links in our predicted network are more likely to have been observed with viruses (Figure 2C; Figure SI4). It is possible that species with more links in the global viral sharing network are more important for viral sharing, and thus have been more likely to be observed with a (zoonotic) virus. Species that are more central in our predicted network could therefore be prioritised for zoonotic surveillance or sampling in the event of viral outbreaks with unknown mammalian origins. Given that mammal diversity predicts patterns of livestock disease³² and zoonoses²², the geographic patterns of network centrality predicted here (Figure 3, Figure SI9; see below) could also be used as a coarse predictor of disease risk to livestock and human health. Similarly, where there is limited knowledge of mammalian host range for newly-discovered viruses, our modelled network can be used to prioritise the sampling of additional species for viral discovery.

It is possible that the high predicted centrality of known hosts was due partly to selective sampling (i.e., viral researchers are more likely to sample wide-ranging and common host

species that also share viruses with many other species^{10,20}). This possibility is supported by the greater link number for species that appear in both EID2 and our dataset rather than in only one of the two, as these species are presumably more well-known (Figure 2C). Similarly, while we believe that our model was successful at accounting for variation in host-level diversity and study effort that influences network topology (see above; Figure S11), there are certain inherent biases in the training data which must be considered when interpreting our findings. Most notably, sharing estimates in our dataset may be affected by the fact that zoonotic discovery efforts commonly search limited geographic regions for a specific virus or group of viruses, artificially increasing the likelihood of detecting these viruses in the same region compared to a geographically random sampling regime. Moreover, when a mammal species (e.g., a bat) is found with a focal virus (e.g., an ebolavirus), it is logical for researchers to then investigate similar, closely related species in nearby locales. These sampling approaches could disproportionately weight the network towards finding phylogeographic associations. However, it is highly encouraging that our model predicted patterns in the external EID2 dataset, which was constructed using different data compilation methods but also comprising a large global body of data covering several decades of research³¹. We believe that our approach is a conservative method for minimising these biases. The knowledge that the observed virome is biased ultimately calls for more uniform sampling across the mammal class and increased coverage of rarely-sampled groups, lending support to ongoing efforts to systematically catalogue mammalian viral diversity³.

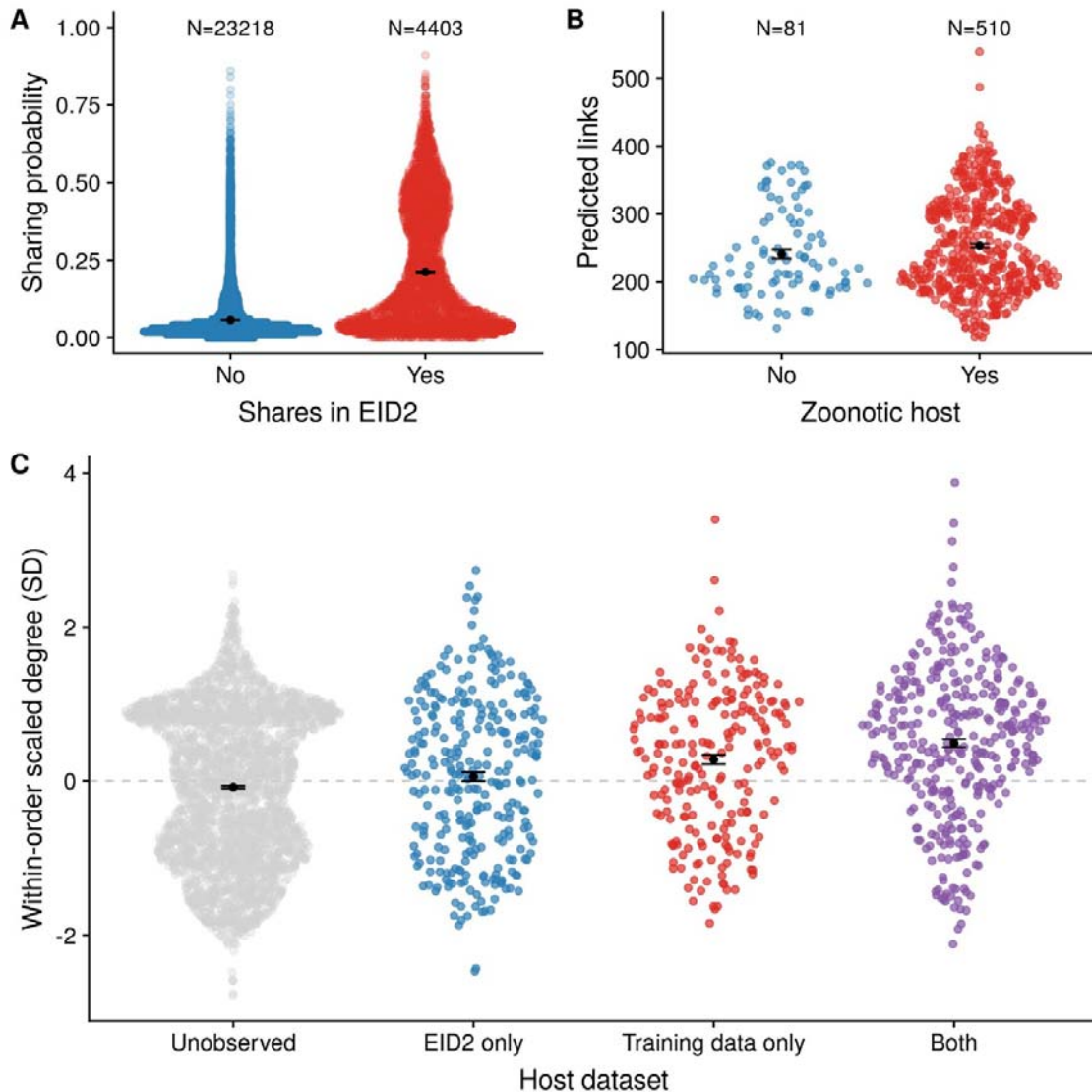


Figure 2: The modelled mammalian viral sharing network predicts observed viral sharing trends. In all figures, points are jittered along the x axis according to a density function; the black points and associated error bars are means \pm standard errors. A: species pairs with higher predicted viral sharing probability from our model were more likely to be observed sharing a virus in the independent EID2 dataset. This comparison excludes species pairs that were also present in our training data. B: species that hosted a zoonotic virus in our dataset had more viral sharing links in the predicted all-mammal network than those without zoonotic viruses. C: species that had never been observed with a virus have fewer links in the predicted network than species that hosted viruses in the EID2 dataset only, in our training data only, or in both. The y axis represents viral sharing link number, scaled to have a mean of 0 and a standard deviation of 1 within each order for clarity. Figure S14 displays these same data without the within-order scaling.

Taxonomic and geographic patterns of predicted viral sharing

Our network predicted strong taxonomic patterns of viral connectedness. 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). In agreement with previous findings, these patterns demonstrate that bats and rodents are among the most important taxa for viral sharing, with ungulates and carnivores somewhat less important⁵. Here both orders' high centrality emerged purely as a result of their phylogenetic diversity and geographic distributions, rather than from other species-level phenotypic traits such as behaviour⁵, life history⁹, or metabolic idiosyncracies³³. 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 size³⁴), this may contribute to the high viral diversity found in bats and rodents⁵. Efforts to prioritise viral sampling regimes should consider biogeography and mammal-mammal interactions in addition to searching for species-level traits associated with high viral diversity.

Partitioning the network into within- and between-order sharing links revealed differences in taxonomic and geographic patterns of viral sharing across taxonomic scales. Bats' and rodents' large numbers of within-order links are driven by large order size (Figure 3C). Interestingly, when these 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. The disparity in within- versus out-of-order link number was particularly large for rodents. This difference 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 SI9). Encouragingly, our network showed predictable scaling laws similar to those of other known ecological networks³⁵. Link numbers in within-order subnetworks (e.g., between different bat species) correlated strongly with species diversity within each order ($R^2 \sim 0.85$), following a power law with a Z value of ~ 0.8 (Figure SI5). Similarly, out-of-order links (e.g., between a bat and a rodent) scaled with an offset 1:1 relationship with the product of the size of both orders (Figure SI6).

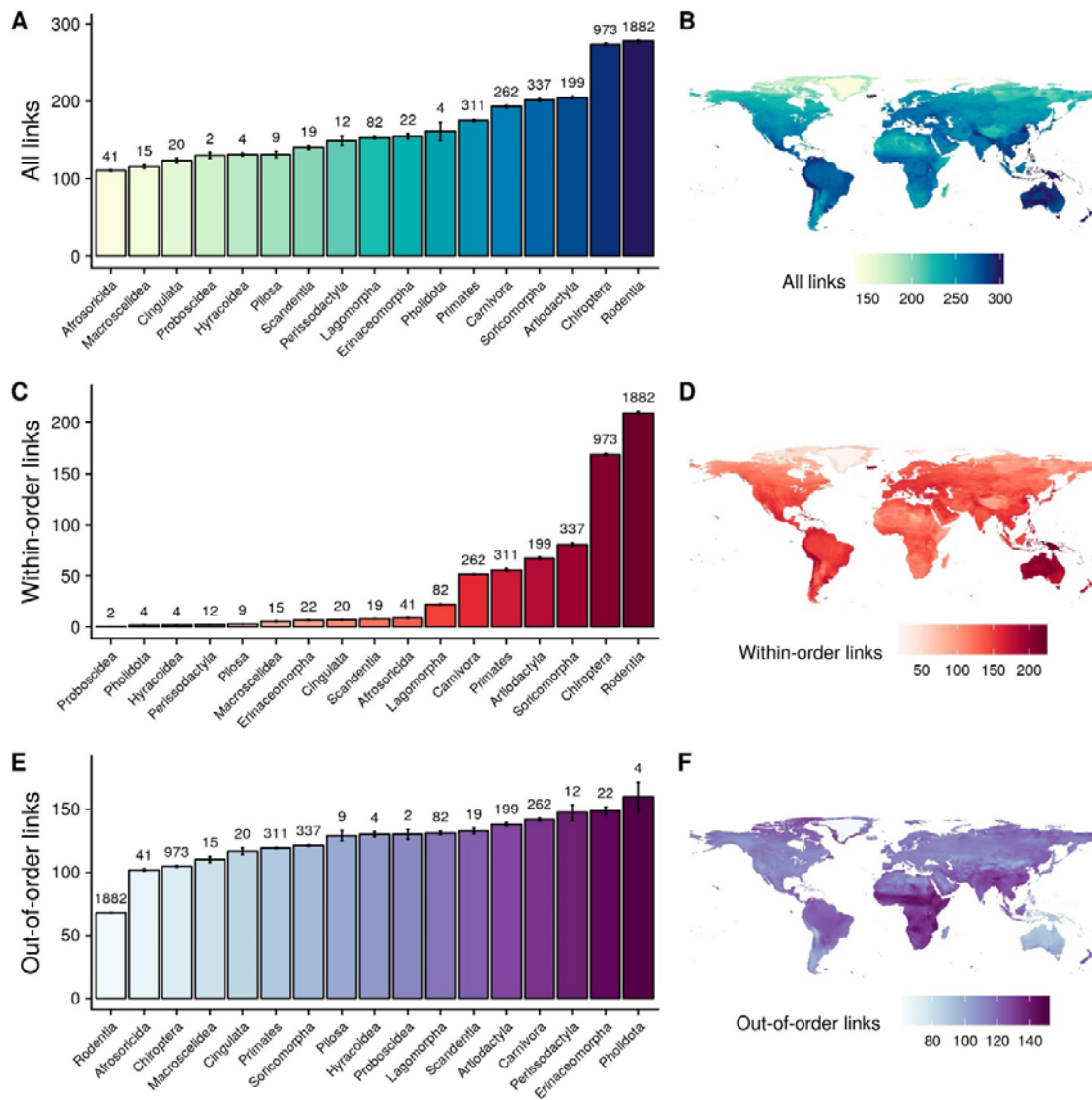


Figure 3: Taxonomic and geographic patterns of mean predicted viral link numbers. Top row: all links; middle row: links with species in the same order; bottom row: links with species in another order. A,C,E: average species-level link numbers for mammalian orders in our dataset. Bars represent means; error bars represent standard errors. B,D,F: geographic distributions of mean link numbers. Distributions were derived by summing the link numbers of all species inhabiting a 25km² grid square and dividing them by the number of species inhabiting the grid square, giving mean degree number at the grid level.

To visualize geographic patterns of viral sharing, we projected species-level link numbers across the species' ranges, summarizing to grid cell-level mean link number (Figure 3B), as well as summed link number (Figure SI9). Average connectedness peaked in tropical areas of

South and Central America, Sub-Saharan Africa, and Southeast Asia, especially in the Andes and Himalayas (Figure 3B). These patterns align with previously-reported hotspots of emerging zoonoses and predicted viral diversity^{5,22} and imply that areas of high biodiversity are centres of viral sharing not just because of the number of overlapping species (i.e., high species richness), but also because the network is more highly connected in these areas (Figure SI9). This densely-connected network structure and the increased biomass present in the tropics might have synergistic implications for cross-species maintenance and transmission of viral diversity in these areas. The geographic distributions of mean predicted within- and between-order viral links differed notably from the distribution of interspecific links generally: the relative importance of South America and East Asia was higher for within-order links (Figure 3D), while Sub-Saharan Africa remained a hotspot for out-of-order links (Figure 3F). Geographic patterns of summed link numbers more closely mirrored underlying host species richness, whether for all links, within-order links, or out-of-order links (Figure SI9).

We acknowledge that our phylogeographic model of viral sharing does not account for complex ecological interactions such as apparent competition, coinfection, or coevolution, all of which will impact how patterns of exposure and host susceptibility translate to realised viral diversity. Future investigations may extend our framework to simulate the dynamic co-speciation of mammals and their viruses in order to account for these processes and/or to explicitly investigate how viral sharing connectivity and viral diversity are correlated across mammal species. Our model may also prove useful for building and parameterising much-needed multi-host network models for conservation purposes, particularly where there is scarce prior information on interspecific pathogen sharing^{34,36}.

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 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 (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 our phylogeographic model required only $\frac{1}{4}$ as many sampling targets in order to identify the correct sharing host. The GAMM 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 dataset³¹: 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.

We observed substantial variation in our model's ability to predict known hosts among viral species. For example, the correct host was predicted first in every iteration for 7 viruses and in the top 10 hosts for 42 viruses. Results for 128 viruses had the focal host falling within the top 100 guesses, and for only 6 viruses were the model-based host searches worse than chance (focal host ranked lower than 50% of all mammals in terms of sharing probability). We used this measure of “predictability” to investigate whether certain viral traits affected the ease with which phylogeography predicted their hosts. Viruses with broad host phylogenetic ranges, most notably Ebola virus, challenge reservoir prediction efforts since

many more species must often be sampled before identifying the correct host(s). To investigate whether the predictive strength of our model was limited for viruses with broad host ranges and/or other viral traits, we fitted a linear mixed model (LMM) which showed a strong negative association between viruses' known phylogenetic host breadth and the predictability of focal hosts (model $R^2=0.70$; host breadth $R^2=0.67$; Figure SI8). This association demonstrates, unsurprisingly, that predicting the hosts of generalist viruses is intrinsically difficult for our model-based method. This adds a potential limitation to the applicability of our network approach, given that zoonotic viruses commonly exhibit wide host ranges^{2,5}. A family-level random effect accounted for little of the apparent variance in predictability among viral families (Figure SI7).

Once viral host range was accounted for, hosts of vector-borne viruses were slightly easier to predict than non-vector-borne viruses ($R^2=0.1$; Figure SI8) – perhaps because the sharing of vector-borne viruses depends more heavily on host geographic distributions (Figure SI3). Despite additional variation in the data, no other viral traits (e.g., RNA vs. DNA, segmented vs. non-segmented) were important in the LMM. This implies that phylogeographic traits are a good broad-scale indicator of viral sharing, particularly when ecological specifics are unknown.

Concluding remarks

In summary, we present a simple, highly interpretable model that predicted a substantial proportion of viral sharing across mammals and is capable of identifying species-level sampling priorities for viral surveillance and discovery. It is also worth noting that the analytical framework and validation we describe here was conducted on a global scale, while many zoonotic sampling efforts occur on a national or regional scale. Restricting the focal mammals to a regional pool may improve the applicability of our model in certain sampling contexts, and future studies could leverage higher-resolution phylogenetic and geographic data to fine-tune predictions. In particular, the mammalian supertree³⁰ has relatively poor resolution at the species tips such that relatedness estimates based on alternative molecular evidence (e.g., full host genome data) may allow more precise estimates of the phylogenetic relatedness effect on viral sharing. Alternatively, our model could be augmented with additional host, virus, and pairwise traits, using similar pairwise formulations of viral sharing

as a response variable, to identify ecological specificities that are critical for the transmission of certain viruses, to partition viral subtypes, and, ultimately, to increase the accuracy of host prediction. By generalising the spatial and phylogenetic processes that drive viral sharing, our model serves as a useful guide for the prioritization of viral sampling, presenting a baseline for future modelling efforts to compare against and improve upon.

Our ability to model and predict macroecological patterns of viral sharing is important in an era of rapid global change. Under all conceivable global change scenarios, many mammals will shift their geographic ranges, whether of their own volition or through human assistance. Mammalian parasite communities will likely undergo considerable rearrangement as a result, with potentially far-reaching ecological consequences³⁷⁻³⁹. Our findings suggest that novel species encounters will provide opportunities for interspecific viral transmission, which could be facilitated by even relatively small changes in range overlap (Figure 1B). Such cross-species viral transmission could have profound implications for conservation and public health, potentially devastating populations of host species without evolved resistance to the novel viruses (e.g., red squirrel declines brought about by parapoxvirus infections spread by introduced grey squirrels⁴⁰) or increasing zoonotic disease risk by introducing viruses to human-adjacent amplifier hosts (e.g., horses increasing the risk of human infection with Hendra virus²⁰). Thus, our global model of mammalian viral sharing provides a crucial complement to ongoing work modelling the spread of hosts, vectors, and their associated diseases as the result of climate change-induced range expansions^{27,37}.

Acknowledgements

This work was conducted during a placement funded by the National Environmental Research Council (NERC) Overseas Research Fund awarded to GFA. GFA's PhD studentship was likewise funded by NERC (grant number: NE/L002558/1). EAE, NR, and KJO were funded by the generous support of the American people through the United States Agency for International Development (USAID) Emerging Pandemic Threats PREDICT project. Additional support was provided by the National Institute of Allergy and Infectious Diseases of the National Institutes of Health (Award Number R01AI110964) and the US Department of Defense, Defense Threat Reduction Agency (HDTRA11710064). We thank Colin Carlson, Verity Hill, and members of EcoHealth Alliance for advice on viral macroecology and for helpful comments on early versions of this work.

References

1. Woolhouse, M. E. J. & Gowtage-Sequeria, S. Host range and emerging and reemerging pathogens. *Emerg. Infect. Dis.* **11**, 1842–7 (2005).
2. Johnson, C. K. *et al.* Spillover and pandemic properties of zoonotic viruses with high host plasticity. *Sci. Rep.* **5**, 1–8 (2015).
3. Carroll, D. *et al.* The Global Virome Project. *Science*. **359**, 872–874 (2018).
4. Carlson, C. J., Zipfel, C. M., Garnier, R. & Bansal, S. Global estimates of mammalian viral diversity accounting for host sharing. *Nat. Ecol. Evol.* **3**, 1070–1075 (2019).
5. Olival, K. J. *et al.* Host and viral traits predict zoonotic spillover from mammals. *Nature* **546**, 646–650 (2017).
6. Han, B. A. *et al.* Undiscovered Bat Hosts of Filoviruses. *PLoS Negl. Trop. Dis.* **10**, e0004815 (2016).
7. Babayan, S. A., Orton, R. J. & Streicker, D. G. Predicting reservoir hosts and arthropod vectors from evolutionary signatures in RNA virus genomes. *Science*. **362**, 577–580 (2018).
8. Luis, A. D. *et al.* A comparison of bats and rodents as reservoirs of zoonotic viruses: are bats special? *Proc. R. Soc. B Biol. Sci.* **280**, 20122753 (2013).
9. Han, B. A., Schmidt, J. P., Bowden, S. E. & Drake, J. M. Rodent reservoirs of future zoonotic diseases. *Proc. Natl. Acad. Sci.* **112**, 7039–7044 (2015).
10. Plourde, B. T. *et al.* Are disease reservoirs special? Taxonomic and life history characteristics. *PLoS One* **12**, e0180716 (2017).
11. Dallas, T. A. *et al.* Host traits associated with species roles in parasite sharing networks. *Oikos* **128**, 23–32 (2019).
12. Plowright, R. K. *et al.* Pathways to zoonotic spillover. *Nat. Rev. Microbiol.* **15**, 502–510 (2017).
13. Ge, X. Y. *et al.* Isolation and characterization of a bat SARS-like coronavirus that uses the ACE2 receptor. *Nature* **503**, 535–538 (2013).
14. Huang, S., Bininda-Emonds, O. R. P., Stephens, P. R., Gittleman, J. L. & Altizer, S. Phylogenetically related and ecologically similar carnivores harbour similar parasite assemblages. *J. Anim. Ecol.* **83**, 671–680 (2014).
15. Wells, K. *et al.* Global spread of helminth parasites at the human-domestic animal-wildlife interface. *Glob. Chang. Biol.* **24**, 3254–3265 (2018).
16. Stephens, P. R. *et al.* Parasite sharing in wild ungulates and their predators: Effects of phylogeny, range overlap, and trophic links. *J. Anim. Ecol.* **88**, 1017–1028 (2019).
17. Streicker, D. G. *et al.* Host Phylogeny Constrains Cross-Species Emergence and Establishment of Rabies Virus in Bats. *Science*. **329**, 676–679 (2010).
18. Willoughby, A. R., Phelps, K. L., Predict Consortium, predict@ucdavis.edu & Olival, K. J. A

- comparative analysis of viral richness and viral sharing in cave-roosting bats. *Diversity* **9**, 1–16 (2017).
19. Davies, T. J. & Pedersen, A. B. Phylogeny and geography predict pathogen community similarity in wild primates and humans. *Proc. R. Soc. B Biol. Sci.* **275**, 1695–1701 (2008).
 20. Glennon, E. E. *et al.* Domesticated animals as hosts of henipaviruses and filoviruses: A systematic review. *Vet. J.* **233**, 25–34 (2018).
 21. Chua, K. B. *et al.* Nipah Virus: A Recently Emergent Deadly Paramyxovirus. *Science*. **288**, 1432–1435 (2000).
 22. Allen, T. *et al.* Global hotspots and correlates of emerging zoonotic diseases. *Nat. Commun.* **8**, 1124 (2017).
 23. Gómez, J. M., Nunn, C. L. & Verdú, M. Centrality in primate-parasite networks reveals the potential for the transmission of emerging infectious diseases to humans. *Proc. Natl. Acad. Sci.* **110**, 7738–41 (2013).
 24. Guy, C., Thiagavel, J., Mideo, N. & Ratcliffe, J. M. Phylogeny matters: Revisiting ‘a comparison of bats and rodents as reservoirs of zoonotic viruses’. *R. Soc. Open Sci.* **6**, 181182 (2019).
 25. Pybus, O. G., Tatem, A. J. & Lemey, P. Virus evolution and transmission in an ever more connected world. *Proc. R. Soc. B Biol. Sci.* **282**, 20142878 (2015).
 26. Sanjuán, R., Nebot, M. R., Chirico, N., Mansky, L. M. & Belshaw, R. Viral mutation rates. *J. Virol.* **84**, 9733–9748 (2010).
 27. Ryan, S. J., Carlson, C. J., Mordecai, E. A. & Johnson, L. R. Global expansion and redistribution of Aedes-borne virus transmission risk with climate change. *PLoS Negl. Trop. Dis.* **13**, e0007213 (2019).
 28. Drake, J. M. & Beier, J. C. Ecological niche and potential distribution of *Anopheles arabiensis* in Africa in 2050. *Malar. J.* **13**, 213 (2014).
 29. IUCN. The IUCN Red List of Threatened Species. *IUCN Red List of Threatened Species. Version 2019-2* (2019). Available at: <https://www.iucnredlist.org>.
 30. Fritz, S. A., Bininda-Emonds, O. R. P. & Purvis, A. Geographical variation in predictors of mammalian extinction risk: big is bad, but only in the tropics. *Ecol. Lett.* **12**, 538–549 (2009).
 31. Wardeh, M., Risley, C., Mcintyre, M. K., Setzkorn, C. & Baylis, M. Database of host-pathogen and related species interactions, and their global distribution. *Sci. Data* **2**, 150049 (2015).
 32. Wang, Y. X. G. G. *et al.* Phylogenetic structure of wildlife assemblages shapes patterns of infectious livestock diseases in Africa. *Funct. Ecol.* **33**, 1332–1341 (2019).
 33. Xie, J. *et al.* Dampened STING-Dependent Interferon Activation in Bats. *Cell Host Microbe* **23**, 297–301 (2018).
 34. Lloyd-smith, J. O. *et al.* Epidemic Dynamics at the Human-Animal Interface. *Science*. **326**, 1362–1368 (2009).

35. Brose, U., Ostling, A., Harrison, K. & Martinez, N. D. Unified spatial scaling of species and their trophic interactions. *Nature* **108**, 167–171 (2003).
36. Silk, M. *et al.* Integrating social behaviour, demography and disease dynamics in network models: applications to disease management in declining wildlife populations. *Philos. Trans. R. Soc. B* **374**, 20180211 (2019).
37. Carlson, C. J. *et al.* Parasite biodiversity faces extinction and redistribution in a changing climate. *Sci. Adv.* **3**, e1602422 (2017).
38. Williams, J. E. & Blois, J. L. Range shifts in response to past and future climate change: Can climate velocities and species' dispersal capabilities explain variation in mammalian range shifts? *J. Biogeogr.* **45**, 2175–2189 (2018).
39. Chen, I.-C., Hill, J. K., Ohlemüller, R., Roy, D. B. & Thomas, C. D. Rapid range shifts of species associated with high levels of climate warming. *Science*. **333**, 1024–6 (2011).
40. Tompkins, D. M., Sainsbury, A. W., Nettleton, P., Buxton, D. & Gurnell, J. Parapoxvirus causes a deleterious disease in red squirrels associated with UK population declines. *Proc. R. Soc. B Biol. Sci.* **269**, 529–533 (2002).
41. R Core Team. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. (2018).
42. Wood, S. N. Fast stable restricted maximum likelihood and marginal likelihood estimation of semiparametric generalized linear models. *J. R. Stat. Soc. Ser. B (Statistical Methodol.* **73**, 3–36 (2011).

Methods

Making the training data network

Code and data for all analyses is available at

<https://github.com/gfalbery/ViralSharingPhylogeography>. Our dataset included 1920 mammal-virus associations obtained from an exhaustive literature search which has been used to investigate how species traits influence mammalian viral diversity⁵. We removed humans and rabies virus from the dataset as both were disproportionately well-connected, and we removed 20 non-Eutherian mammals because they were extreme phylogenetic outliers, leaving 591 Eutherian mammals that shared 401 viruses. We made an unweighted bipartite network using the mammal-virus associations and projected the unipartite mammal-mammal network, which we then converted into a sequence of all unique mammal-mammal pairs where 1/0 denoted whether the pair of species shared a virus or not. This comprised only the lower triangle of the adjacency matrix to avoid duplicating associations and to remove self-

connections, and only included mammals with at least one sharing link (final N=174345 unique mammal-mammal pairs). 6.4% of these pairs shared at least one virus.

All analyses were performed in R version 3.6.0⁴¹. Phylogenetic similarity was calculated using a mammalian supertree³⁰ as previously described⁵. Pairwise phylogenetic distances were defined as the cumulative branch length between the two species and were scaled to between 0 and 1, and subtracted from 1 to give a measure of relative phylogenetic similarity (rather than distance). Of the 4716 Eutherian species in the mammalian supertree, 591 had virus association records in our fully-connected network and 4196 had known geographic ranges. We used IUCN species ranges to quantify species' geographic distributions²⁹. These range maps are generated based on expert knowledge and only comprise species presence/absence information rather than density. We converted all range polygons to 25 KM² raster grids. For each species-pair, we quantified range overlap as the number of raster grid squares jointly inhabited by the two species (in the Mollweide projection, which exhibits equal grid size), divided by the total number of grid squares occupied by these species combined, so that each value was scaled from 0-1: $\text{overlap}_{A,B} = \text{grid}_{A,B} / (\text{grid}_A + \text{grid}_B - \text{grid}_{A,B})$. Research effort was quantified as previously described, using species' disease-related citation number⁵. To fit citation number as a pairwise trait, we took the smaller of a pair of species' respective citations, and log-transformed the value. Domestication status was defined *sensu lato*, again as previously described⁵, based on whether a species was ever seen in a domestic setting. We fit this as a binary pairwise trait where 1=at least one of the species was domesticated and 0=neither species had been domesticated.

Estimating factors driving viral sharing with (GAMM)

We fitted a Generalised Additive Mixed Model (GAMM) to examine which traits lead to viral sharing among mammal pairs using accelerated discretized implementation in the **mgcv** package⁴². We fitted viral sharing (0/1) as the response variable, with a binomial family specification. The model had the following structure:

$$\begin{aligned} \text{Bernoulli}(\text{Viral sharing}) \sim & s(\text{Phylogenetic similarity, by} = \text{ordered}(\text{Gz})) + \\ & t2(\text{Phylogenetic similarity, Geographic overlap, by} = \text{ordered}(!\text{Gz})) + \\ & \text{Minimum citation number} + \text{Domestication status} + \\ & \text{mm} (\text{Species 1} + \text{Species 2}) \end{aligned}$$

The first term (“s”) represents a phylogeny effect smooth fitted across species pairs that did not overlap in space ($Gz=1$), and “t2” represents a phylogeny:geography tensor product smooth fitted to species that had geographic overlap greater than zero ($Gz=0$). This allowed us to model these two aspects of the data separately. “mm” represents a multi-membership random effect, accounting for the identity of both species in the pair. We implemented this multi-membership effect to control for species-level effects by including a species-level effect for both the row (Species 1) and column (Species 2) of the sharing matrix. Using the paraPen specification in **mgcv**, these random effects were constrained to sample from the same distribution, resulting in a single estimate of the variance associated with each unique species. Most precisely, these species-level effects in our model help capture variation in viral sharing that could likely be explained by factors that are unobserved or otherwise excluded (i.e., species-level differences in underlying viral diversity, which would be expected to positively impact the probability of interspecific sharing). In sum, this model formulation allowed us to estimate the effect of pairwise predictors (geographic overlap, phylogenetic similarity) in determining viral sharing as well as evaluate the influence of species identity.

To investigate whether the effects of geography and phylogeny depended on which subset of viruses we investigated, we repeated our GAMM analysis on non-exclusive subnetworks of mammal-mammal pairs based on the types of viruses they were connected by. Viral subtypes included RNA viruses, vector-borne RNA viruses, non-vector-borne RNA viruses, and DNA viruses. There were only 2 vector-borne DNA viruses, limiting our ability to test this group as a separate dataset. We eliminated from each analysis any hosts that were not carrying the focal virus type.

Validating the GAMM

To check the fit of the model, we predicted 0/1 viral sharing values from the model 1000 times and examined how the values compared to the proportions of 0's and 1's in the observed data, finding high agreement between the two. We repeated this procedure using a) the full dataset; b) only the fixed effects, with random effects randomised in each iteration; and c) only the random effects, with fixed effects held at the mean. We then used these predicted links to create 1000 unipartite viral sharing networks, estimating link numbers

(degree centrality) for species in each replicated network. We took the mean of these values across the 1000 replicated networks to give the predicted values displayed in Figure S11.

We quantified deviance contributions of our explanatory variables by calculating model deviance when dropping each variable, and comparing these against full model and an intercept-only model deviance. For each of our explanatory variables (geographic overlap, phylogenetic similarity, minimum citation number, domestication status, and species-level random effects) we randomised the observed values 1000 times, then predicted sharing probabilities for these values using our model estimates. This randomisation procedure allowed us to predict while accounting for the uneven data distribution, rather than using mean values. We subtracted the full model deviance from the calculated deviances and then calculated the proportions of each to obtain the percentage deviance explained by each variable. We divided the full model deviance by that estimated from the intercept-only model to calculate the deviance explained by the full GAMM.

Simulating viral sharing networks

Following reconstruction of the observed network as part of our model validation, we repeated the prediction process on our 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 species²⁹ 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 biodiversity.

We validated the predicted network by comparing it to sharing patterns in the Enhanced Infectious Diseases Database (EID2)³¹. 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 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 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 viral hosts, we iteratively investigated the sharing patterns of known hosts independently for all viruses with >1 host. We removed one host at a time, and then investigated which species the remaining 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 viruses with the remaining species, our model was taken to be useful for predicting patterns of mammal sharing based on known host distributions. The mean ranking of the focal hosts across each prediction iteration was used as a measure of “predictability” for each virus. We carried out this process for the 250 viruses with more than one known host with associated geographic and phylogenetic data and then on the 109 such viruses in the EID2 data.

Once the predictability of each virus was calculated, we fitted a linear mixed model examining $\log_{10}(\text{mean focal host rank})$ as an inverse measure of predictability (higher rank corresponds to decreased predictability) for each virus. We added mean phylogenetic host similarity as a fixed effect and viral family as a random effect to quantify how viral phylogeny affected predictability. We included additional viral traits in the model, including: cytoplasmic replication (0/1); segmentation (0/1); vector-borne transmission (0/1); double- or single-strandedness; DNA or RNA; enveloped or non-enveloped; or zoonotic ability (0/1 for whether the virus was associated with humans in our dataset).

Supplementary Information

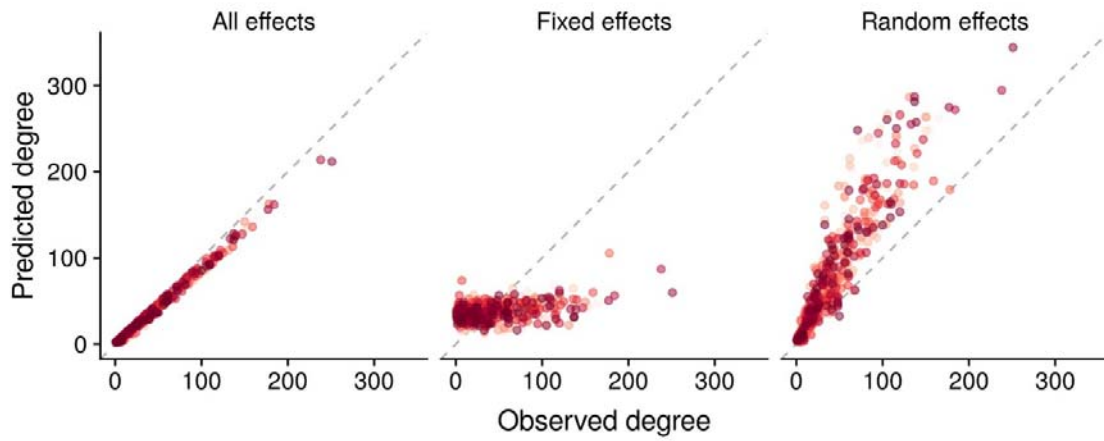


Figure S11: Predicted degree centrality of species in our training data network, predicted using our GAMM estimates. Fixed + random effects were very effective at reproducing individual species' degree centrality (left); fixed effects were less effective (middle); and random effects alone had a strong but imperfect effect (right).

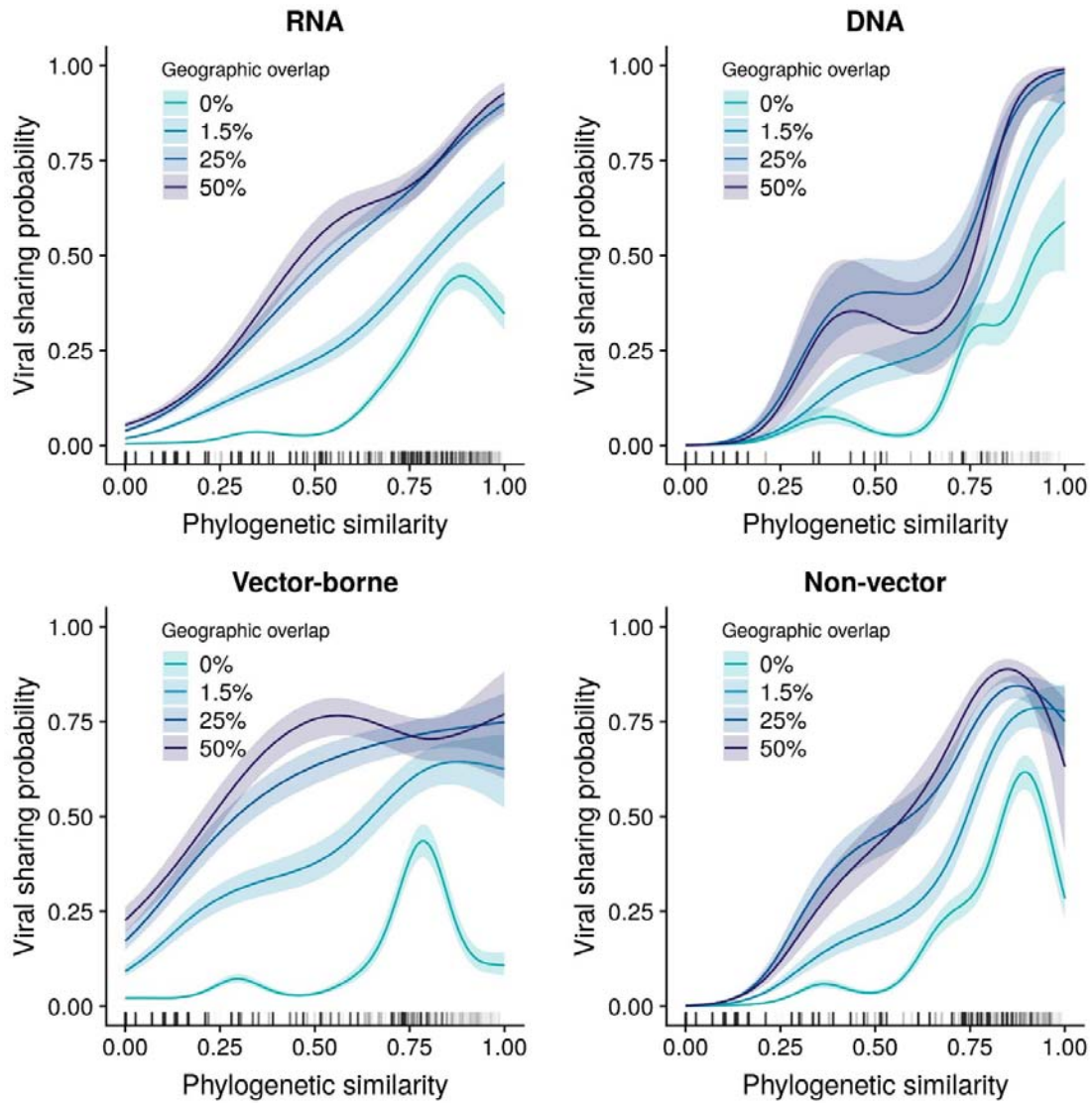


Figure S12: GAMM-derived viral sharing estimates for the effect of phylogenetic similarity for four viral subsets (top row: all RNA viruses and DNA viruses; bottom row: vector-borne RNA viruses and non-vector-borne RNA viruses). Each GAMM smooth is displayed at multiple geographic overlap values (different colours).

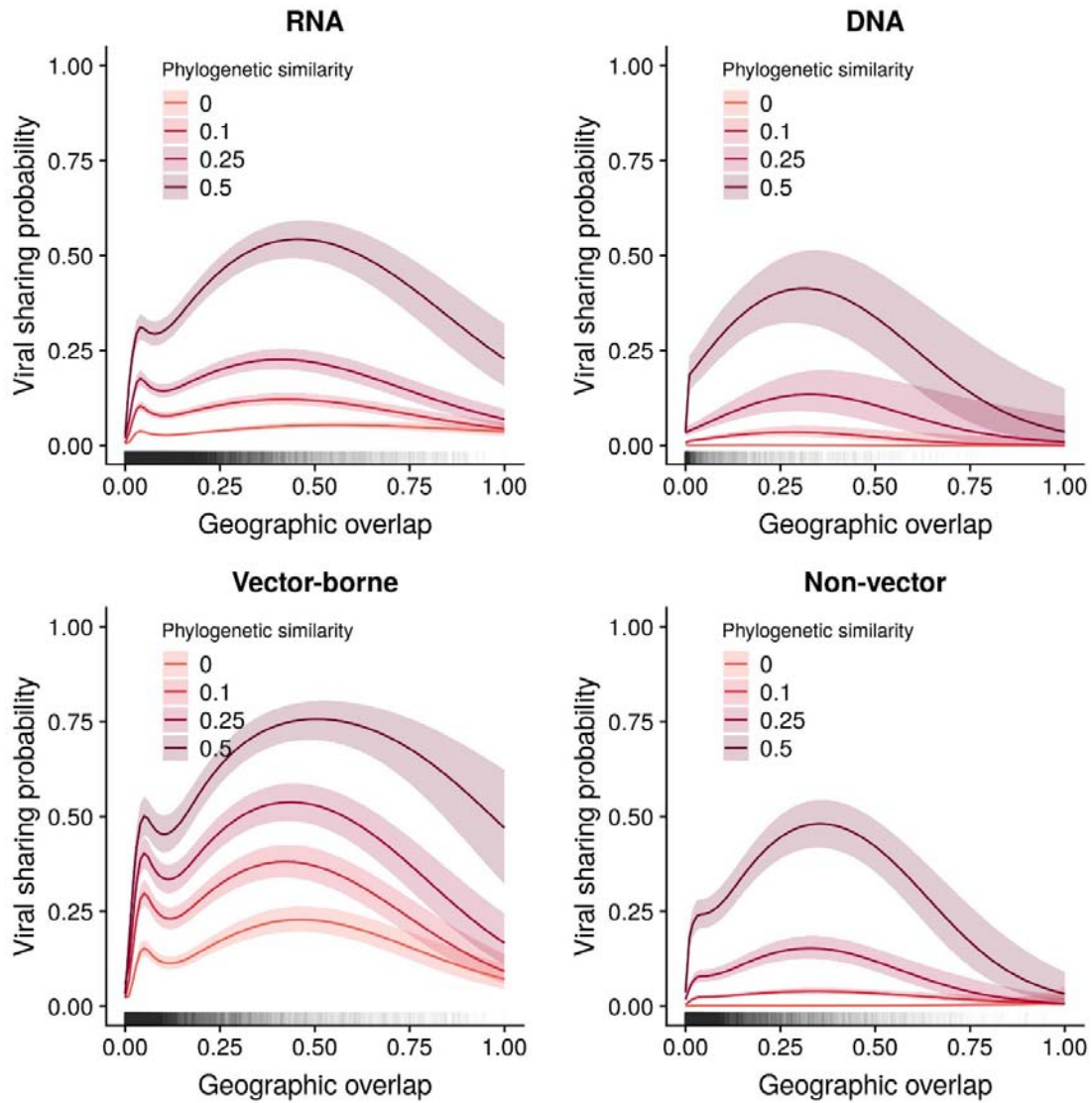


Figure SI3: GAMM-derived viral sharing estimates for the effect of geographic overlap for four viral subsets (top row: all RNA viruses and DNA viruses; bottom row: vector-borne RNA viruses and non-vector-borne RNA viruses). Each GAMM smooth is displayed at multiple geographic overlap values (different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

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

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

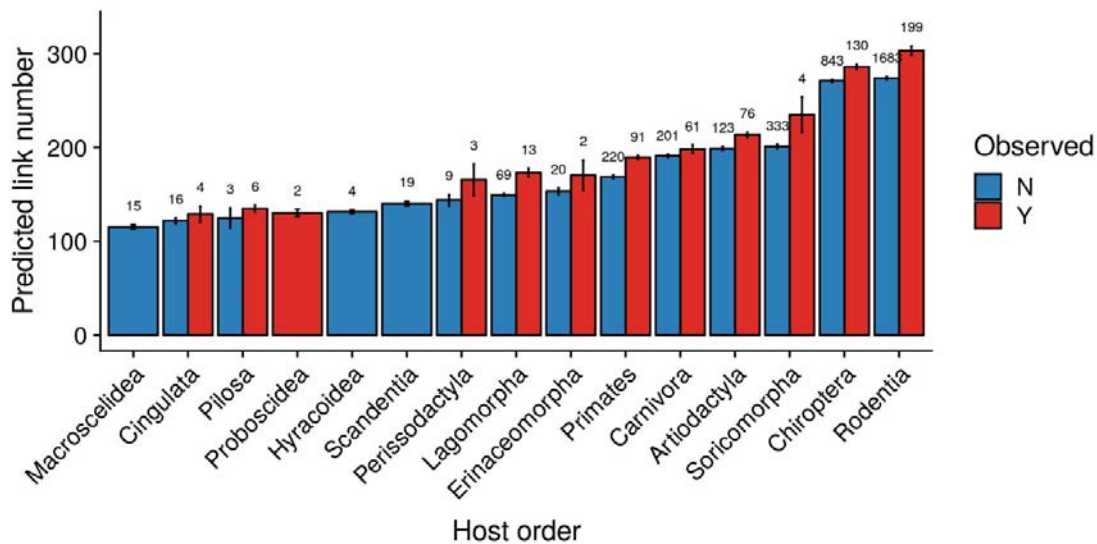


Figure SI4: Mammal species that were observed with at least one virus in the training dataset or the EID2 dataset had higher degree centrality (link number) in our predicted network. This figure displays the raw data that are displayed in Figure 2C in the main text, but without being scaled within orders.

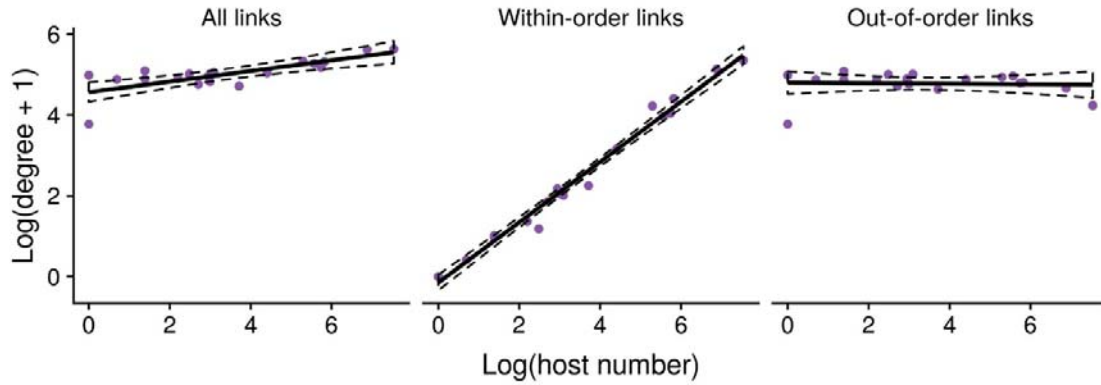


Figure SI5: Scaling of degree centrality (link numbers) followed a power law when looking within-orders, but not between orders. The trend line and 95% confidence intervals are derived from a linear model fitted to the data.

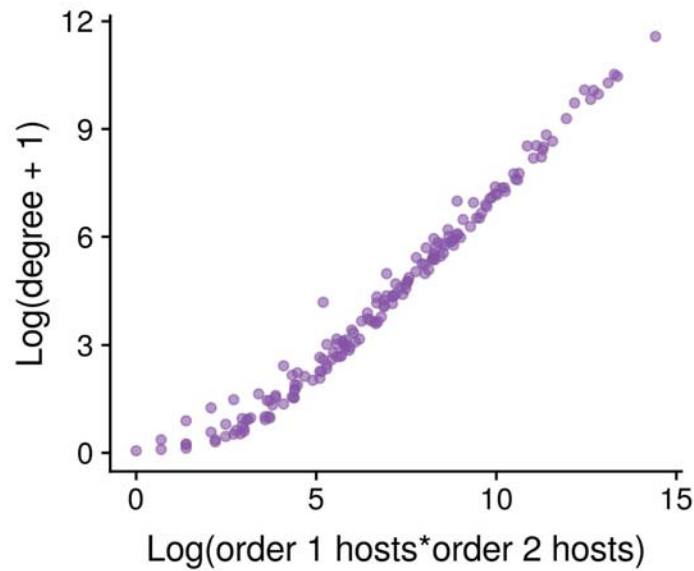


Figure SI6: Predicted between-order link numbers scales according to the log-product of the number of species in the two orders. Each point represents a pair of orders (N=171).

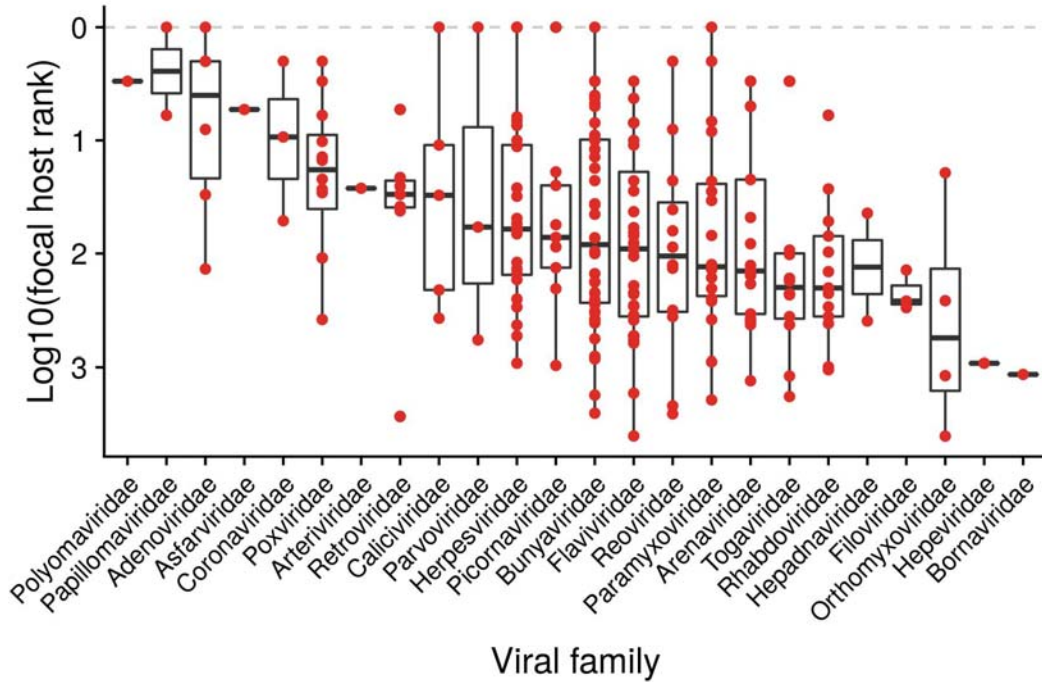


Figure SI7: The phylogeographic predictability of viruses' reservoir hosts varied considerably across viral families, although the family-level random effect did not account for much of the model's variance. Families are ordered along the x axis in order of decreasing predictability. The y axis displays the mean rank of the focal host in our reservoir host prediction simulation, on a reversed log10-scale. Values closer to the top of the figure represent more predictable viruses.

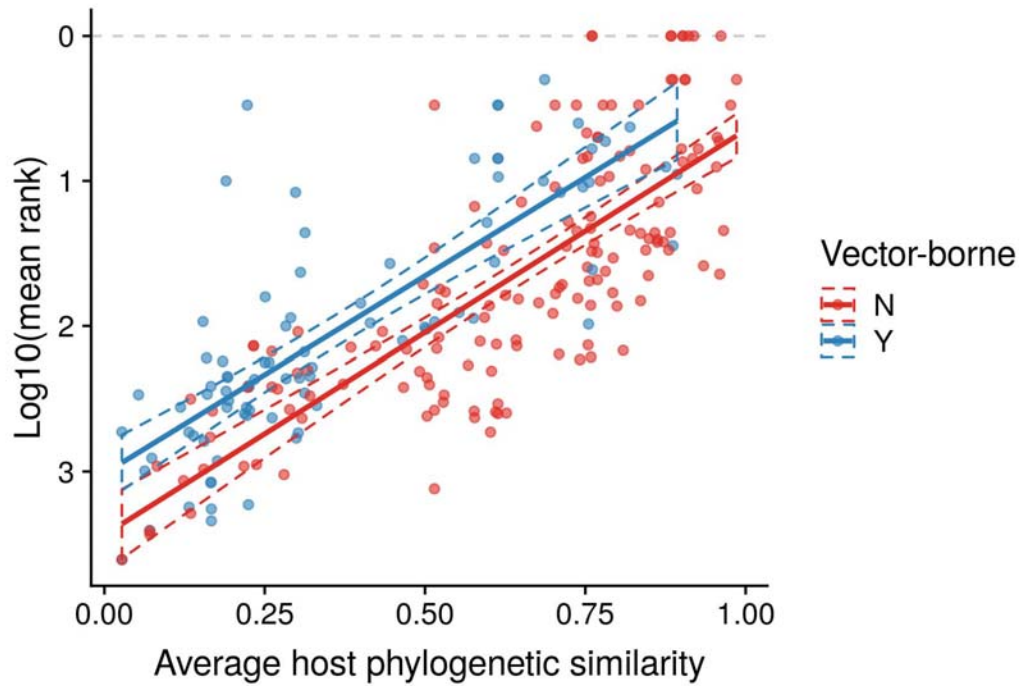


Figure S18: Viral host range strongly impacted the predictability of reservoir hosts. The x axis displays the mean phylogenetic similarity of a virus's hosts (i.e., an inverse measurement of viral host range). The y axis displays the mean rank of the focal host in our reservoir host prediction simulation, on a reversed log₁₀-scale. Values closer to the top of the figure represent more predictable viruses. The trend lines and 95% confidence intervals were derived from a linear mixed model fitted to the data.

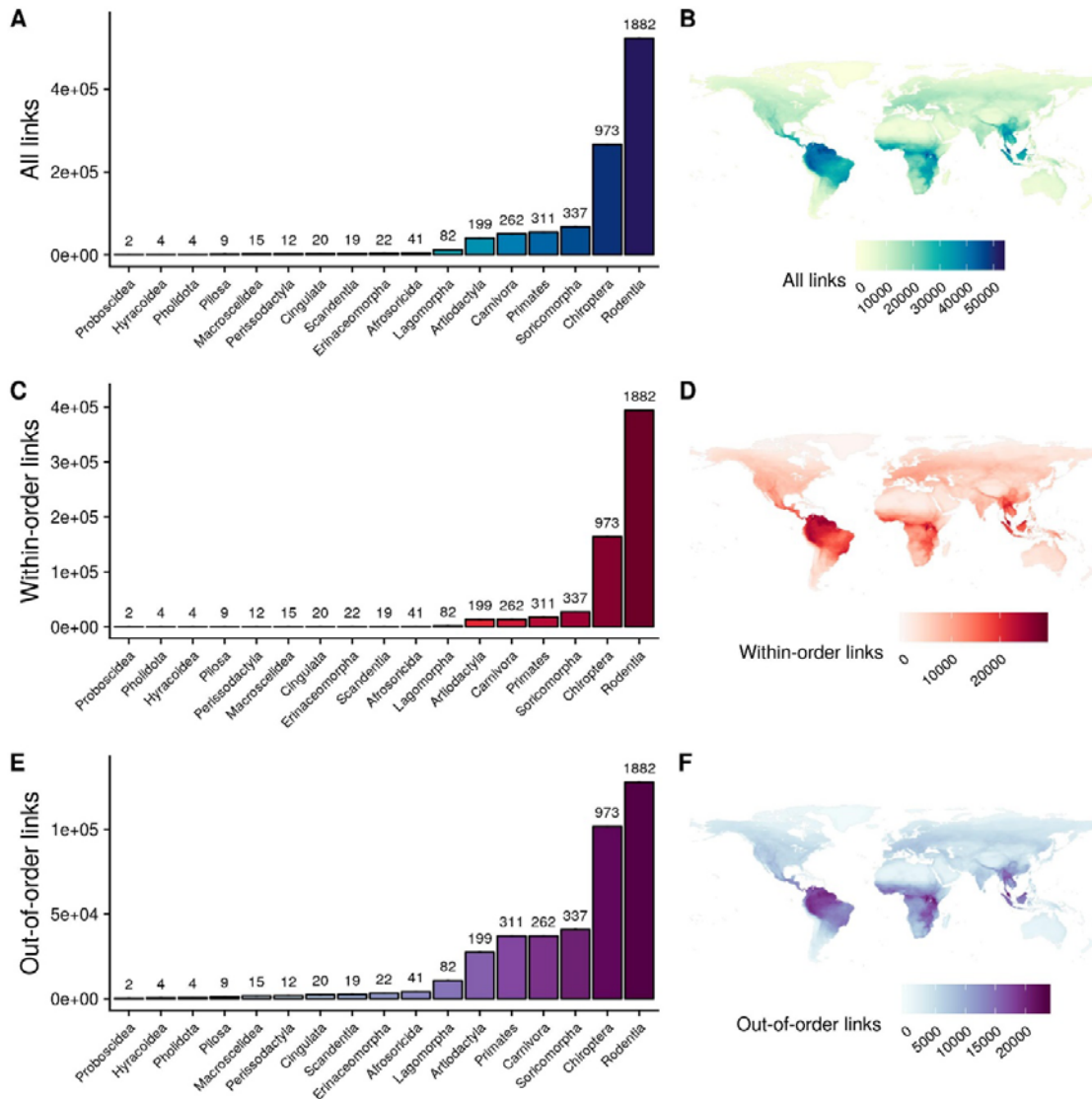


Figure SI9: Taxonomic and geographic patterns of predicted viral link numbers. Top row: all links; middle row: links with species in the same order; bottom row: links with species in another order. A,C,E: summed species-level link numbers for mammalian orders in our dataset. B,D,F: geographic distributions of link numbers. Distributions were derived by summing the link numbers of all species inhabiting a grid square.