

Predicting the tripartite network of mosquito-borne disease

Tad A. Dallas^{a,b,†}, Sadie J. Ryan^{c,d}, Ben Bellekom^e, Anna Fagre^f, Rebecca Christofferson^a and Colin J. Carlson^{g,h}

^a*Department of Biological Sciences, Louisiana State University, Baton Rouge, LA 70802*

^b*Department of Biological Sciences, University of South Carolina, Columbia, SC, 29208*

^c*Department of Geography, University of Florida, Gainesville, FL 32611 USA*

^d*Emerging Pathogens Institute, University of Florida, Gainesville, FL 32611, USA*

^e*Department of Zoology, 11a Mansfield Road, Oxford OX1 3SZ, UK*

^f*Department of Microbiology, Immunology, and Pathology, Colorado State University, Fort Collins, CO 80523, US.A.*

^g*Center for Global Health Science and Security, Georgetown University Medical Center, Washington, D.C. 20007*

^h*Department of Microbiology and Immunology, Georgetown University Medical Center, Washington, D.C. 20007*

[†]*Corresponding author: tad.a.dallas@gmail.com*

Running title: A tripartite model of vector-borne disease

Author contributions: TAD, CJC, and SJR conceived the study. TAD and CJC performed the analysis. All authors contributed to manuscript writing.

Acknowledgements: Lauren Holian provided comments and discussion on an earlier version of this work. This work has been supported by funding to the Viral Emergence Research Initiative (VERENA) consortium, including a grant from the U.S. National Science Foundation (NSF BII 2021909).

Conflict of interest: The authors have no conflicts of interest to declare.

Keywords: Arbovirus, Host-virus, Link prediction, Networks, Vector-borne

1

2 **Predicting the tripartite network of mosquito-borne disease**

3 **Abstract**

4 The potential for a pathogen to infect a host is mediated by traits of both the host
5 and pathogen, as well as the complex interactions between them. Arthropod-borne
6 viruses (arboviruses) require an intermediate arthropod vector, which introduces
7 an additional layer of compatibility filters. Existing computational models for the
8 prediction of host-virus networks rarely incorporate the unique aspects of vector
9 transmission, instead treating vector biology as a hidden, unobserved layer. Here,
10 we explore two possible extensions to existing approaches, to address this nuance:
11 first, we added vector traits into predictions of the bipartite host-virus network;
12 and second, we used host, vector, and virus traits to predict the tripartite host-
13 vector-virus network. We tested both approaches on the most thoroughly charac-
14 terized group of arboviruses; mosquito-borne flaviviruses of mammals, including
15 dengue, yellow fever, and Zika virus. Using host-virus models, we find that the
16 inclusion of vector traits may improve inference in some cases, while viral traits
17 proved to be the most important for model performance. Further, we found that
18 it was possible to predict full life cycles (host-vector-virus links), but the model
19 only showed fair performance, and was heavily influenced by the geographic bias of
20 component input datasets (especially the dipteran biting data). Both approaches
21 are interesting avenues for further model development, but our results keenly un-
22 derscore a need to collect more comprehensive datasets to characterize arbovirus
23 ecology, across a wide geographic scope, especially outside of North America, and
24 to better identify molecular traits that underpin host-vector-virus interactions.

25

26 Introduction

27 Emerging viruses continue to pose a threat to human and wildlife populations [1].
28 A growing set of computational tools have explored viral dynamics in the context
29 of species interaction networks using a set of tools called *link prediction models*.
30 Typically, these represent hosts and viruses as a bipartite network of either known
31 interactions (that occur in nature [2, 3]) or all possible interactions (including, for
32 example, experimental infections [4]), with both represented as links in the network
33 [5]. Host-virus link prediction models are predominantly trained on the genomic,
34 immunological, morphological, and ecological traits of hosts and viruses (e.g., [6,
35 7]), while some approaches also leverage information on the latent structure of
36 the network instead of, or in addition to, these traits [8, 9]. The objective of
37 these modeling exercises is to learn about the underlying biology, explain and
38 reproduce patterns found in nature, and anticipate what future dynamics of viral
39 emergence could look like. For example, many models use networks to understand
40 why some viruses can infect humans but others cannot, with the objective of
41 identifying animal viruses that could someday infect humans for the first time.
42 In most cases, these models assume that any given “link” between a host and a
43 virus could represent a self-contained transmission cycle (though not necessarily
44 *onwards* transmission, e.g., West Nile virus in humans and horses [10]).

45 Vector-borne disease (VBD) transmission substantially complicates this con-
46 ceptual framework. Vector-borne viruses require an additional species—usually an
47 arthropod (hence arthropod-borne viruses, or *arboviruses*)—to move them between
48 hosts, which adds complexity into their ecology, epidemiology, and evolution. For
49 example, in the case of arboviruses, the presence of both virus and suitable hosts
50 is not necessarily sufficient for transmission, and the presence or absence of suit-
51 able vectors (e.g., their geographic distributions or host preferences) may be a
52 latent variable in ecological datasets [11]. Moreover, the “compatibility filters”
53 that can be inferred from the host-virus network will be incomplete, as models
54 will miss both the molecular and physiological determinants of vector-virus com-
55 patibility (i.e., vector competence) and the behavioral and ecological determinants

56 of vector-host compatibility (i.e., biting preferences, in the case of blood-feeding
57 arthropods). If vectors are entirely omitted from the inference process, a model
58 might therefore reach spurious conclusions about whether a given host and virus
59 are incompatible based on their biology, or otherwise miss key drivers of network
60 structure; for example, arboviruses have been shown repeatedly to have a higher-
61 than-expected host breadth [12].

62 No one canonical approach exists to address vector transmission in link predic-
63 tion studies. Vector transmission could be described as a binary trait of viruses,
64 which may help make some distinctions (e.g., separating the ecology of mosquito-
65 borne and tick-borne flaviviruses from counterparts like hepatitis C), but leaves
66 much to be desired in terms of information content (e.g., not distinguishing the
67 tick- and mosquito-borne flaviviruses). The possibility of incorporating more de-
68 tailed information on vector-borne transmission into these models has been under-
69 explored, likely because arboviruses are usually seen as a complicated exception to
70 existing datasets, rather than a feature with significant impacts on network struc-
71 ture. Incorporating traits characterizing the life cycle of arboviruses might improve
72 model performance, given that virus traits are often sparser than host traits, and
73 their interactions usually have non-additive but positive effects on model perfor-
74 mance. However, adding sparse traits that only describe some of the viruses in the
75 network could also reduce accuracy if the network includes a mix of vector-borne
76 and directly-transmitted viruses.

77 Alternately, vectors could be added directly into the network as an additional
78 layer of nodes (Figure 1). While previous work has predicted vector-virus networks
79 [13], none have predicted *host-vector-virus* networks. Existing network models
80 have been used to predict undetected links in *tripartite* networks [14], but this has
81 yet to be explored for ecological networks. This approach would be much more
82 informative than the bipartite form, but also requires difficult-to-obtain data: syl-
83 vatic VBD cycles tend to be characterized one at a time in scientific literature (e.g.,
84 “*Culex quinquefasciatus* vectors West Nile virus in house finches”). While available
85 datasets could be used to reconstruct these cycles from each of their component

86 parts (biting preferences, vector competence, and host-virus compatibility), to our
87 knowledge, this has not previously been explored in predictive work.

88 To address this, we developed two new approaches and tested them on mosquito-
89 borne flaviviruses, a well-studied group that includes important zoonoses like
90 dengue, West Nile, yellow fever, and Zika viruses. Through a synthesis of ex-
91 isting data sources, we combined data on mammal-virus associations [12], vector-
92 flavivirus associations [13], and diptera-mammal biting preferences [15]. We com-
93 bined these data into one mammal-mosquito-flavivirus network, which can also be
94 reduced down to a mammal-flavivirus network where viruses' mosquito commu-
95 nities are represented as node metadata. Using boosted regression trees (BRT;
96 a machine learning method popular in ecological modeling, also sometimes called
97 gradient boosting machines), we tested two approaches to predicting vector-borne
98 transmission as an aspect of the host-virus network. First, we predicted the
99 mammal-flavivirus network using every possible combination of host, vector, and
100 virus traits, as metadata for any given host-virus association, assuming that ad-
101 ditional data layers would enhance model performance. This was generally shown
102 to be true, although the combination of host and vector trait data was not infor-
103 mative compared to the incorporation of viral trait data. Second, we developed a
104 tripartite model of vector-borne disease transmission, in which each link represents
105 a known host-vector-virus link and attempted to predict those complete cycles us-
106 ing traits of hosts, mosquito vectors, and viruses. We found that these models
107 performed more poorly on average, but that they were able to make better than
108 random predictions, including some of relevance to arboviral ecology and human
109 health.

110 **Methods**

111 **Host, vector, and virus data** Host-virus interaction data were obtained from
112 the CLOVER database [16], a manually- and programmatically-curated database
113 of host-virus associations built by reconciling four disparate datasets (the Host-

114 Parasite Phylogeny Project, or HP3 [12]; the Global Mammal Parasite Database
115 v2.0 [17]; the Enhanced Infectious Disease Database [18]; and an unnamed dataset
116 curated by Shaw *et al.* [19]). We used CLOVER release 0.1.2, which includes
117 data on 5,477 known interactions between 831 viruses of 1,085 mammal species.
118 These data have been carefully cleaned for taxonomic quality control and include
119 detailed metadata on interaction evidence. These data are also part of a larger
120 open database called The Global Virome in One Network (VIRION), the largest
121 open atlas of vertebrate-virus associations [20]. Although more data is available
122 from this source, we restricted our analysis to the manually-curated data to prevent
123 inclusion of spurious interactions.

124 Vector-virus association data were taken from a previous study that aimed to
125 predict the mosquito-flavivirus network. [13] These data include 334 associations
126 between 180 mosquito species and 37 flaviviruses. Host-vector association data
127 were taken from a recent study of dipteran biting networks [15]. These data
128 describe 1744 associations between 255 biting dipteran species and 214 hosts (in-
129 cluding 67 mammals). Trait data for hosts, vectors, and viruses were assembled
130 from published sources. Thirty-three traits on mosquito life history, ecology, and
131 geography and 22 traits on viral features, were taken from the Evans *et al.* study
132 of the mosquito-flavivirus network [13]. Finally, we used a total of 18 traits on
133 mammal life history, ecology, and morphology from the PanTHERIA database
134 [21].

135 **Modeling approach** Boosted regression tree (BRT) models were used to model
136 host-virus and host-vector-virus associations. BRT models have previously been
137 used to model species distributions [22], predict associations in bipartite networks
138 [23, 24, 25, 5], and in other conservation and management settings e.g., [26]. Much
139 of the diversity of applications can be attributed in part to the allowance for nonlin-
140 ear responses and variable interactions in BRT models. Since the regression tree
141 is hierarchical, “upstream” splits based on one variable influence “downstream”
142 splits, which automatically models variable interactions. Further, the process of
143 boosting enhances learning on complex data, as the process produces many regres-

144 sion trees with a small number of splits, each of these “weak learners” iteratively
145 build on previous trees to account for the remaining variation. This approach
146 removes the need to partition variance among submodels, as the goal is not to
147 examine the components of variance explained, but to assess overall model per-
148 formance with the inclusion or exclusion of particular variable sets. Models were
149 trained in the *R* statistical programming language [27] using the *gbm* package [28].

150 **Model 1: Modeling mammal-virus associations as a bipartite network**

151 We used the mammal and virus trait data as described above. However, mosquito
152 vector “traits” were created by calculating the number of mosquito species in a
153 given genus which were demonstrated to transmit a particular flavivirus [13]. This
154 is because each host-virus association could be transmitted by any number of
155 mosquito species, creating a range of trait values that may be less informative
156 than simply knowing breadth and composition of the vector community. This
157 resulted in a total of 19 mosquito vector covariates, ranging in value from 0 to 22
158 species. We removed covariates with less than 25% data coverage, resulting in 13
159 host traits, 19 mosquito covariates (as virus traits), and 17 virus traits.

160 The data were split into 80% training and 20% testing sets, where model per-
161 formance was assessed on the 20% test set. A total of 20 models per covariate
162 group were fit in order to account for the random train/test split. These same 20
163 train/test divisions were used across the different covariate models, as we trained
164 every possible combination of host, vector, and virus trait data to predict host-
165 virus associations. Together, this resulted in a dataset that allows the estimation
166 of the relative influence of host traits, viral traits, and vector community data on
167 resulting mammal-virus associations. We sampled background data by randomly
168 combining host and virus species, resulting in 25% known positive associations
169 and 75% background data.

170 We subset these data in two different ways, to explore how vector data may
171 improve prediction of 1) flaviviruses for which we have some vector data (235
172 known host-virus associations) and 2) all vector-borne viruses (3016 host-virus

173 associations). This breakdown corresponds to data subsets of 1) only mosquito-
174 borne flaviviruses present in [13] and 2) all viruses that were recorded as vector-
175 borne (or unknown) in the Clover data [16]. We present the flavivirus-specific
176 results here, which are qualitatively similar to the more general models for all
177 vector-borne viruses, which are in the Supplemental Materials.

178 **Model 2: Modeling mammal-mosquito-virus associations as a tripartite**
179 **network** Using the same data resource as used above on host-virus associations,
180 we now considered the identity of the mosquito vector species, and the association
181 between the vector and virus [13], and the feeding association between mosquito
182 vector and mammal species [15]. While host and virus traits were largely the same
183 as considered above, the mosquito vector traits consisted of a set of 33 mosquito
184 vector traits from [13]. Host and virus traits must have 75% of data coverage –
185 the same as in *Model 1* – to be included in this analysis. This resulted in 8 host
186 traits, 29 vector traits, and 16 virus traits. A tripartite link – detailing the full
187 host-vector-virus cycle – was only considered if there were all three associations;
188 host-vector association, vector-virus association, and host-virus association. This
189 creates a situation where a host and vector species may interact, and that vector
190 may be infected by a virus, but this is not a confirmed link if there is no evidence
191 that the host is infected by the virus.

192 A total of 135 full tripartite links were documented. We sampled background
193 data by randomly combining host, vector, and virus species and then adding
194 enough unique host-vector-virus background points to have 50% true tripartite
195 links and 50% background data. Models were trained in the same manner as in
196 **Model 1**.

197 **Assessing model performance** Model performance was quantified using two
198 measures; accuracy and the area under the receiver operating characteristic (AUC).
199 Accuracy was defined as the correctly estimated positives (true positives) and neg-
200 atives (true negatives) over all the predictions, capturing the fraction of times the
201 model correctly classified host-virus associations in the holdout data. Accuracy is

202 bounded between 0 and 1, where larger values correspond to higher model perfor-
203 mance. AUC is a widely used metric of model discrimination that captures the
204 ability of the classifier to rank positive instances higher than negative instances.
205 AUC is bounded between 0 and 1, where a random model will perform with AUC
206 of 0.5 on average, and values closer to 1 indicate higher model performance.

207 **Data and code availability** *R* code and data to reproduce the analyses is
208 available on figshare at
209 <https://doi.org/10.6084/m9.figshare.17033309>.

210 Results

211 **Model 1: The mammal-virus models** Models trained only on host (AUC
212 $= 0.57$) or vector ($AUC = 0.46$) traits consistently performed poorly at the task
213 of host-virus link prediction (Figure 2), though the viral trait model performed
214 well ($AUC = 0.95$). Generally, combinations of predictor features led to improved
215 model performance. The full model including host, vector, and virus traits per-
216 formed extremely well ($AUC = 0.96$). However, both the host-virus and vector-
217 virus traits only models also performed extremely well (performance differences
218 among these models were essentially indistinguishable; Figure 2). The inclusion of
219 viral traits seems to have been particularly important; for comparison, the model
220 using host and vector traits to predict host-virus associations barely performed
221 better than random ($AUC = 0.59$).

222 Variables important for predicting host-virus associations were generally con-
223 served across submodels considering all combinations of host, vector, and virus
224 traits (Figure 3). In the full bipartite model, the most informative variable was
225 whether a virus was found in the Pacific region (likely a proxy for Zika virus,
226 which spread through Pacific islands preceding the epidemic in the Americas).
227 Other important characteristics predictive of host-virus associations in bipartite
228 models including virus traits were disease severity, genome length, year of virus

229 isolation, if the virus is found in Africa or Australia, and viral clade. In models
230 that omitted virus traits, the top predictors represented host allometry (body mass
231 and metabolic rate, an unsurprising axis of variation) and *Culex* association, which
232 likely captures a latent split between some bird-reservoired viruses (e.g., West Nile
233 virus) and primate-reservoired ones (e.g., dengue and Zika virus).

234 Overall, our results suggest that models learned from vector trait data, par-
235 ticularly in the full model, where the contribution of each individual variable is
236 more diffuse. However, our findings also indicate that the inclusion of vector data
237 only minimally improved performance after data on hosts and viruses was already
238 available. As host-virus models are usually trained only on host and virus trait
239 data, our findings suggest that the incorporation of vector data into a host-virus
240 model is an imperfect way to explore the role of vectors in structuring the host-
241 virus network. However, this also suggests that improved arthropod trait data
242 could improve model performance, and thus the importance of the vector cannot
243 be overlooked.

244 Finally, we investigated whether including vector trait data would improve per-
245 formance even if only available for a subset of data informing the network. To test
246 this, we trained the model on a network that included all the arboviruses present
247 in the CLOVER dataset, even though viral trait data and vector associations were
248 only known for flaviviruses. We found that the model using just host and virus
249 traits performed substantially worse here ($AUC = 0.70$) than the flavivirus-only
250 model with those traits ($AUC = 0.95$). We found that the best performing models
251 were those that used vector and virus traits ($AUC = 0.98$) and those that included
252 host, vector, and virus traits ($AUC = 0.99$; Figure 2). We suggest that this finding
253 indicates that adding data on the vector aspect of transmission may be useful even
254 when it only covers a subset of species in the network.

255 **Model 2: The tripartite model** Models trained on tripartite (i.e., host-vector-
256 virus) associations had moderate explanatory power (mean $AUC = 0.64$ (0.065);
257 mean Accuracy = 0.66 (0.046) out of 100 models trained on random subsets). This

258 lower model performance could simply be due to the smaller amount of data used
259 for training (recall that only 135 full tripartite links were known), or the imbalance
260 between the number of potential full tripartite links given host, vector, and virus
261 diversity, and the small number of realized links (see the small number of red links
262 in Figure 4). Although the model’s performance was only fair, we found that the
263 model still predicted higher suitability for tripartite links where one or two of the
264 three possible components were confirmed (Figure 5), even though these would
265 be recorded as a “0” outcome variable the same as if none of them were known.
266 We suggest that this indicates the model was identifying and reproducing real
267 biological signals of compatibility.

268 The top nine covariates to predicting tripartite (i.e., host-vector-virus) associ-
269 ations were host ($n = 5$) or virus ($n = 4$) traits (Figure 6). The top predictors
270 mostly reflected the geography of transmission (host geographic range size, virus
271 transmission in Asia, vector presence in Africa), the life history of the host (age
272 at first birth, lifespan, weaning age, and neonate body mass), and aspects of viral
273 transmission (genome length and transmission by non-mosquito arthropods).

274 The predictions made by the tripartite model suggest the model may be able
275 to recover interesting or important biologically-plausible interactions. Both the
276 top predicted “undiscovered” human-mosquito-virus links (Table 1) and mammal-
277 mosquito-virus links (Table 2) heavily over-represent a small number of viruses,
278 in particular Wesselbron virus and West Nile virus. This is driven by the existing
279 level of sampling in the data: West Nile has the greatest number of known hosts
280 ($n = 103$ species) and mosquito vectors ($n = 51$); Wesselbron has the second high-
281 est number of vectors ($n = 41$), though many fewer hosts ($n = 11$; ranked #13).
282 This “rich-get-richer” has been previously debated as a strength or weakness for
283 link prediction models; it may be that models are identifying a genuine biological
284 signal of generality (which is known to be true for these viruses), but they may
285 also be recapitulating sampling bias [5, 29] and underpredicting link probabilities
286 for undersampled species. Indeed, the richness of flavivirus data available to us in
287 this study is likely largely due to a discovery and data synthesis bump in the wake

288 of the Zika virus epidemic in the Americas. The mammal-mosquito-virus predic-
289 tions also contain a visible signal of geographic bias: most of the top predictions
290 either involve agricultural species (pigs, *Sus scrofa*; cows, *Bos taurus*; or sheep,
291 *Ovis aries*), synanthropic species (black rats, *Rattus rattus*), or charismatic North
292 American species (the opossum, *Didelphis virginiana*; the raccoon, *Procyon lotor*;
293 the white-tailed deer, *Odocoileus virginianus*). These likely reflect a compounded
294 bias between the host-virus association data and the biting data, the latter of
295 which is particularly limited to North American and European species.

296 Despite the signal of data bias in these predictions, the models reveal several
297 predictions of biological interest. For example, *Anopheles hyrcanus* is predicted
298 as a possible vector of Kokobera virus in humans. The virus was implicated in
299 an outbreak of acute polyarticular illness in Australia in the 1980s based on serol-
300 ogy, but it remains poorly understood [30]. The virus was first isolated from
301 *Culex annulirostris*, which also vectors Japanese encephalitis virus and a hand-
302 ful of others; *An. hyrcanus* is a European and Asian mosquito only currently
303 known to vector Japanese encephalitis virus. Similarly, the model predicts that
304 *Culex tritaeniorhynchus* – the main vector of Japanese encephalitis virus, found
305 in southeast Asia – could transmit Murray Valley encephalitis virus in wallabies
306 (*Macropus agilis*). Neither the Australian virus nor the host have been recorded
307 in association with this vector, but as of 2021, the mosquito has been detected in
308 Australia [31], indicating the possibility that this interaction could now emerge.

309 Discussion

310 In this study, we considered two approaches to incorporate arboviral life cycles
311 into link prediction models of the mammal-flavivirus network. First, we used a
312 host-virus (bipartite) framework, and assessed the relative influence of including
313 different trait covariates. We found that viral traits were the strongest contributor
314 to model performance, and the incorporation of host and vector traits into the bi-
315 partite models did little to improve model performance. Second, we explored how

316 these models could be extended to predict the entire host-vector-virus (tripartite)
317 network. This framing is both inherently more complex than the host-virus pre-
318 dictive problem, and is massively limited by the availability of training data, but
319 appears promising for future development.

320 Neither of these approaches provided a complete solution to the host-vector-
321 virus prediction problem, though their limitations differ slightly, with different
322 implications for next steps. Adding vector community data to the bipartite (host-
323 virus) models may be useful where data allow, but may be less important when
324 more detailed, biologically meaningful viral trait data are available. Compared to
325 synthetic datasets of animal ecology, life history, and morphology, only a handful
326 of viral traits (e.g., genome length or disease severity) are available in a standard-
327 ized format, to the point that viral host range is itself often used as a viral trait
328 (e.g., our “primate” or “bird” traits, or “host breadth” (see Table S1)). Recently,
329 some studies have begun to use immunogenetic or genome composition variables
330 to characterize host and virus compatibility more directly [32, 33, 34, 35, 36, 37];
331 comparable features for vectors are not yet available or tested in this framework.
332 Shifting towards these kinds of predictors could help models identify more mean-
333 ingsful signals of virus-animal compatibility, and proportionally reduce the signal
334 of bias in predictions.

335 In contrast, directly modeling the host-vector-virus tripartite network addresses
336 the nuance of vector transmission head-on, but this problem is more severely data
337 limited. As a result, these predictions are very visibly influenced by the geographic
338 and taxonomic bias in the component datasets. However, these data limitations
339 can be addressed by investment in future work characterizing arboviral life cycles in
340 understudied areas [38]. Vector-virus combinations can be tested in the laboratory,
341 including in model-experiment feedback designs that leverage existing predictions
342 (e.g., [13]) much like model-guided fieldwork can be used to optimize viral discovery
343 [25]. Similarly, further investigation of mosquito biting behavior will help resolve
344 the host-vector component [15], highlighting the need for “basic” natural history
345 research even on mosquitoes that are not known to be primary vectors of human

346 disease.

347 Our study is the first to attempt modeling the entire tripartite host-vector-virus
348 network. This is a clear knowledge gap in existing approaches to modeling the
349 host-virus network: identifying a suitable host-pathogen association that has no
350 shared vector may not accurately estimate spillover risk. This may be particularly
351 relevant to efforts to identify viruses with undiscovered zoonotic potential, as the
352 presence or absence of human-biting mosquitoes will be a key contributor to their
353 emergence risk [39]. Similarly, the tripartite framework can provide useful insights
354 into the establishment of sylvatic cycles in interepidemic periods or upon expansion
355 into new geographic areas. The ability of arboviruses to persist in non-human hosts
356 may determine whether an epidemic ends as immunity grows (like Zika virus in
357 the Americas, which was primarily transmitted human-to-human by *Aedes aegypti*
358 and *Ae. albopictus*) or instead becomes a regular occurrence (e.g., yellow fever
359 in the Americas, which is maintained by *Haemogogus spp.* and *Sabethes spp.* in
360 non-human primates, between human epidemics driven by *Aedes aegypti*). These
361 are likely to be particularly important nuances as arboviruses continue to spread
362 around an increasingly globalized world in a changing climate [40, 41, 42, 43]

363 The broader question of “how should we model multi-layer ecological interaction
364 networks” is also one that is likely to have broader implications in computational
365 ecology. For example, there are other cases where researchers are interested the
366 traits that structure tripartite networks, such as bat-bat fly-pathogen networks or
367 plant-pest-parasitoid networks. Multilayer networks are also a topic of increasing
368 interest in network science and mathematics, which will likely open doors for more
369 advanced predictive approaches than the extensions we propose here. This is
370 therefore a promising space for the development of future models, particularly if
371 approached through the lens of iterative validation and data collection [25].

372 References

- 373 [1] Jones, K. E., Patel, N. G., Levy, M. A., Storeygard, A., Balk, D., Gittleman,
374 J. L. & Daszak, P., 2008 Global trends in emerging infectious diseases. *Nature*
375 **451**, 990–993.
- 376 [2] Kading, R. C., Kityo, R. M., Mossel, E. C., Borland, E. M., Nakayiki, T.,
377 Nalikka, B., Nyakarahuka, L., Ledermann, J. P., Panella, N. A., Gilbert,
378 A. T. *et al.*, 2018 Neutralizing antibodies against flaviviruses, babanki virus,
379 and rift valley fever virus in ugandan bats. *Infection Ecology & Epidemiology*
380 **8**, 1439215.
- 381 [3] Fagre, A. C., Lewis, J., Miller, M. R., Mossel, E. C., Lutwama, J. J.,
382 Nyakarahuka, L., Nakayiki, T., Kityo, R., Nalikka, B., Towner, J. S. *et al.*,
383 2021 Subgenomic flavivirus rna (sfrna) associated with asian lineage zika virus
384 identified in three species of ugandan bats (family pteropodidae). *Scientific*
385 *Reports* **11**, 1–8.
- 386 [4] Malmlov, A., Bantle, C., Aboellail, T., Wagner, K., Campbell, C. L., Eckley,
387 M., Chotiwan, N., Gullberg, R. C., Perera, R., Tjalkens, R. *et al.*, 2019
388 Experimental zika virus infection of jamaican fruit bats (*artibeus jamaicensis*)
389 and possible entry of virus into brain via activated microglial cells. *PLoS*
390 *Neglected Tropical Diseases* **13**, e0007071.
- 391 [5] Albery, G. *et al.*, 2021 The science of the host-virus network. *Nature Micro-*
392 *biology* p. in press.
- 393 [6] Dallas, T., Park, A. W. & Drake, J. M., 2017 Predicting cryptic links in
394 host-parasite networks. *PLoS computational biology* **13**, e1005557.
- 395 [7] Elmasri, M., Farrell, M. J., Davies, T. J., Stephens, D. A. *et al.*, 2020 A
396 hierarchical bayesian model for predicting ecological interactions using scaled
397 evolutionary relationships. *Annals of Applied Statistics* **14**, 221–240.

- 398 [8] Ghasemian, A., Hosseinmardi, H., Galstyan, A., Airoidi, E. M. & Clauset, A.,
399 2020 Stacking models for nearly optimal link prediction in complex networks.
400 *Proceedings of the National Academy of Sciences* **117**, 23393–23400.
- 401 [9] Becker, D. J., Albery, G. F., Kessler, M. K., Lunn, T. J., Falvo, C. A., Czirják,
402 G. Á., Martin, L. B. & Plowright, R. K., 2020 Macroimmunology: The drivers
403 and consequences of spatial patterns in wildlife immune defence. *Journal of*
404 *Animal Ecology* **89**, 972–995.
- 405 [10] Bowen, R. A. & Nemeth, N. M., 2007 Experimental infections with west Nile
406 virus. *Current Opinion in Infectious Diseases* **20**, 293–297.
- 407 [11] Huang, Y.-J. S., Higgs, S. & Vanlandingham, D. L., 2019 Arbovirus-mosquito
408 vector-host interactions and the impact on transmission and disease patho-
409 genesis of arboviruses. *Frontiers in microbiology* **10**, 22.
- 410 [12] Olival, K. J., Hosseini, P. R., Zambrana-Torrel, C., Ross, N., Bogich, T. L.
411 & Daszak, P., 2017 Host and viral traits predict zoonotic spillover from mam-
412 mals. *Nature* **546**, 646–650.
- 413 [13] Evans, M. V., Dallas, T. A., Han, B. A., Murdock, C. C. & Drake, J. M., 2017
414 Data-driven identification of potential Zika virus vectors. *Elife* **6**, e22053.
- 415 [14] Chen, C., Grady, S. K., Ellingson, S. R. & Langston, M. A., 2021 Gene-
416 disease-drug link prediction using tripartite graphs. In *Proceedings of the*
417 *12th ACM Conference on Bioinformatics, Computational Biology, and Health*
418 *Informatics*, pp. 1–1.
- 419 [15] Bellekom, B., Hackett, T. D. & Lewis, O. T., 2021 A network perspective on
420 the vectoring of human disease. *Trends in Parasitology* pp. S1471–4922.
- 421 [16] Gibb, R., Albery, G. F., Becker, D. J., Brierley, L., Connor, R., Dallas, T. A.,
422 Eskew, E. A., Farrell, M. J., Rasmussen, A. L., Ryan, S. J. *et al.*, 2021 Data
423 proliferation, reconciliation, and synthesis in viral ecology. *BioScience* p. in
424 press.

- 425 [17] Stephens, P., Pappalardo, P., Huang, S., Byers, J., Farrell, M., Gehman, A.,
426 Ghai, R., Haas, S., Han, B., Park, A. *et al.*, 2017 Global mammal parasite
427 database version 2.0. *Ecology* **98**, 1476–1476.
- 428 [18] Wardeh, M., Risley, C., McIntyre, M. K., Setzkorn, C. & Baylis, M., 2015
429 Database of host-pathogen and related species interactions, and their global
430 distribution. *Scientific data* **2**, 1–11.
- 431 [19] Shaw, L. P., Wang, A. D., Dylus, D., Meier, M., Pogacnik, G., Dessimoz, C.
432 & Balloux, F., 2020 The phylogenetic range of bacterial and viral pathogens
433 of vertebrates. *Molecular ecology* **29**, 3361–3379.
- 434 [20] Carlson, C. J., Gibb, R. J., Albery, G. F., Brierley, L., Connor, R., Dallas,
435 T., Eskew, E. A., Fagre, A. C., Farrell, M. J., Frank, H. K. *et al.*, 2021 The
436 global virome in one network (virion): an atlas of vertebrate-virus associa-
437 tions. *bioRxiv* .
- 438 [21] Jones, K. E., Bielby, J., Cardillo, M., Fritz, S. A., O’Dell, J., Orme, C. D. L.,
439 Safi, K., Sechrest, W., Boakes, E. H., Carbone, C. *et al.*, 2009 PanTHERIA:
440 a species-level database of life history, ecology, and geography of extant and
441 recently extinct mammals. *Ecology* **90**, 2648–2648.
- 442 [22] Elith, J., Leathwick, J. R. & Hastie, T., 2008 A working guide to boosted
443 regression trees. *Journal of Animal Ecology* **77**, 802–813.
- 444 [23] Dallas, T. A., Han, B. A., Nunn, C. L., Park, A. W., Stephens, P. R. &
445 Drake, J. M., 2019 Host traits associated with species roles in parasite sharing
446 networks. *Oikos* **128**, 23–32.
- 447 [24] Dallas, T. A. & Becker, D. J., 2021 Taxonomic resolution affects host- parasite
448 association model performance. *Parasitology* **148**, 584–590.
- 449 [25] Becker, D. J., Albery, G. F., Sjodin, A. R., Poisot, T., Bergner, L. M., Dallas,
450 T. A., Eskew, E. A., Farrell, M. J., Guth, S., Han, B. A. *et al.*, 2021 Optimiz-
451 ing predictive models to prioritize viral discovery in zoonotic reservoirs. *The*
452 *Lancet Microbe* p. in press.

- 453 [26] Soykan, C. U., Eguchi, T., Kohin, S. & Dewar, H., 2014 Prediction of fishing
454 effort distributions using boosted regression trees. *Ecological Applications* **24**,
455 71–83.
- 456 [27] R Core Team, 2020 *R: A Language and Environment for Statistical Comput-*
457 *ing*. R Foundation for Statistical Computing, Vienna, Austria.
- 458 [28] Greenwell, B., Boehmke, B., Cunningham, J. & Developers, G., 2020 *gbm:*
459 *Generalized Boosted Regression Models*. R package version 2.1.8.
- 460 [29] Gibb, R., Albery, G. F., Mollentze, N. F., Eskew, E. A., Brierley, L., Ryan,
461 S. J., Seifert, S. N. & Carlson, C. J., 2021 Mammal virus diversity estimates
462 are unstable due to accelerating discovery effort. *bioRxiv* .
- 463 [30] Boughton, C. R., Hawkes, R. A. & Nairn, H. M., 1986 Illness caused by a
464 kokobera-like virus in south-eastern australia. *Medical journal of Australia*
465 **145**, 90–92.
- 466 [31] Lessard, B. D., Kurucz, N., Rodriguez, J., Carter, J. & Hardy, C. M., 2021
467 Detection of the japanese encephalitis vector mosquito *Culex tritaeniorhynchus*
468 in australia using molecular diagnostics and morphology. *Parasites & vectors*
469 **14**, 1–11.
- 470 [32] Babayan, S. A., Orton, R. J. & Streicker, D. G., 2018 Predicting reser-
471 voir hosts and arthropod vectors from evolutionary signatures in rna virus
472 genomes. *Science* **362**, 577–580.
- 473 [33] Brierley, L. & Fowler, A., 2021 Predicting the animal hosts of coronaviruses
474 from compositional biases of spike protein and whole genome sequences
475 through machine learning. *PLoS Pathogens* **17**, e1009149.
- 476 [34] Mollentze, N., Babayan, S. & Streicker, D., 2021 Identifying and prioritizing
477 potential human-infecting viruses from their genome sequences. *bioRxiv* pp.
478 2020–11.

- 479 [35] Poisot, T., Ouellet, M.-A., Mollentze, N., Farrell, M. J., Becker, D. J., Albery,
480 G. F., Gibb, R. J., Seifert, S. N. & Carlson, C. J., 2021 Imputing the mam-
481 malian virome with linear filtering and singular value decomposition. *arXiv*
482 *preprint arXiv:2105.14973* .
- 483 [36] Wardeh, M., Blagrove, M. S., Sharkey, K. J. & Baylis, M., 2021 Divide-and-
484 conquer: machine-learning integrates mammalian and viral traits with net-
485 work features to predict virus-mammal associations. *Nature Communications*
486 **12**, 1–15.
- 487 [37] Young, F., Rogers, S. & Robertson, D. L., 2020 Predicting host taxonomic
488 information from viral genomes: A comparison of feature representations.
489 *PLoS computational biology* **16**, e1007894.
- 490 [38] Swei, A., Couper, L. I., Coffey, L. L., Kapan, D. & Bennett, S., 2020 Patterns,
491 drivers, and challenges of vector-borne disease emergence. *Vector-Borne and*
492 *Zoonotic Diseases* **20**, 159–170.
- 493 [39] dos Passos Cunha, M., Duarte-Neto, A. N., Pour, S. Z., Ortiz-Baez, A. S.,
494 Černý, J., de Souza Pereira, B. B., Braconi, C. T., Ho, Y.-L., Perondi, B.,
495 Sztajn bok, J. *et al.*, 2019 Origin of the são paulo yellow fever epidemic of
496 2017–2018 revealed through molecular epidemiological analysis of fatal cases.
497 *Scientific reports* **9**, 1–10.
- 498 [40] Ryan, S. J., Carlson, C. J., Mordecai, E. A. & Johnson, L. R., 2019 Global ex-
499 pansion and redistribution of aedes-borne virus transmission risk with climate
500 change. *PLoS neglected tropical diseases* **13**, e0007213.
- 501 [41] Ryan, S. J., Carlson, C. J., Tesla, B., Bonds, M. H., Ngonghala, C. N., Morde-
502 cai, E. A., Johnson, L. R. & Murdock, C. C., 2021 Warming temperatures
503 could expose more than 1.3 billion new people to zika virus risk by 2050.
504 *Global Change Biology* **27**, 84–93.
- 505 [42] Kraemer, M. U., Reiner, R. C., Brady, O. J., Messina, J. P., Gilbert, M.,
506 Pigott, D. M., Yi, D., Johnson, K., Earl, L., Marczak, L. B. *et al.*, 2019 Past

507 and future spread of the arbovirus vectors *Aedes aegypti* and *Aedes albopictus*.
508 *Nature microbiology* **4**, 854–863.

509 [43] Messina, J. P., Brady, O. J., Golding, N., Kraemer, M. U., Wint, G. W.,
510 Ray, S. E., Pigott, D. M., Shearer, F. M., Johnson, K., Earl, L. *et al.*, 2019
511 The current and future global distribution and population at risk of dengue.
512 *Nature microbiology* **4**, 1508–1515.

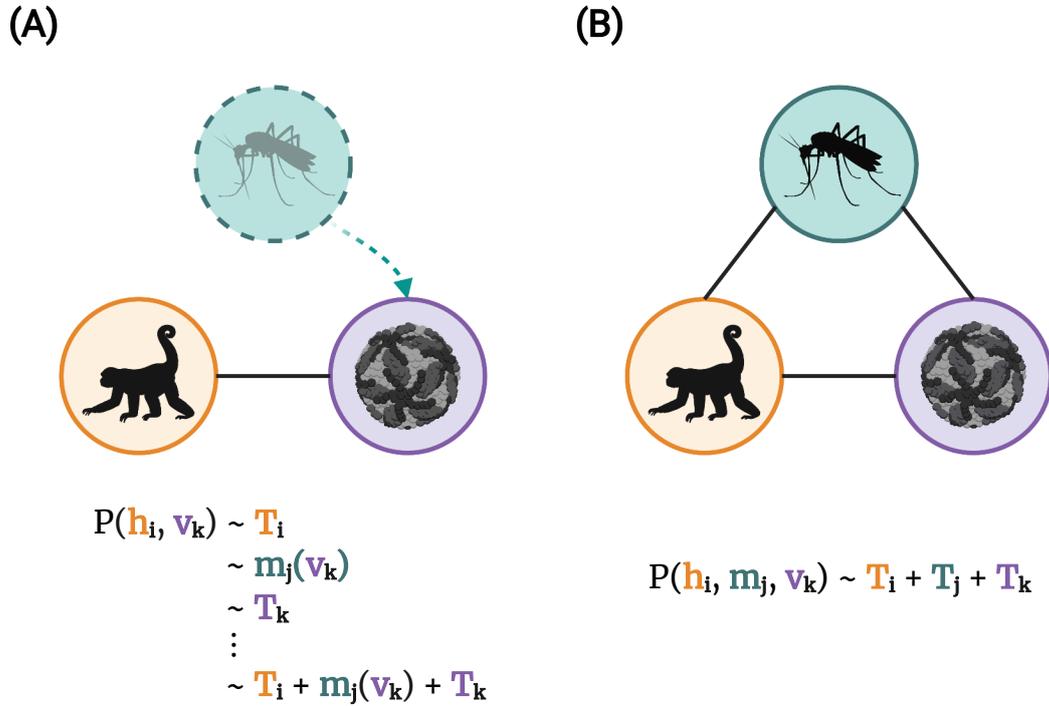


Figure 1: (A) Predicting host-virus associations (a bipartite network) based on host traits (T_h), virus traits (T_v), and vector communities ($m(v)$) associated with viruses, is a different problem than (B) predicting host-vector-virus associations (a tripartite network) based on host traits, vector traits, and virus traits. In this paper, we consider both solutions as approaches to end goals like forecasting potential novel associations or spillover scenarios

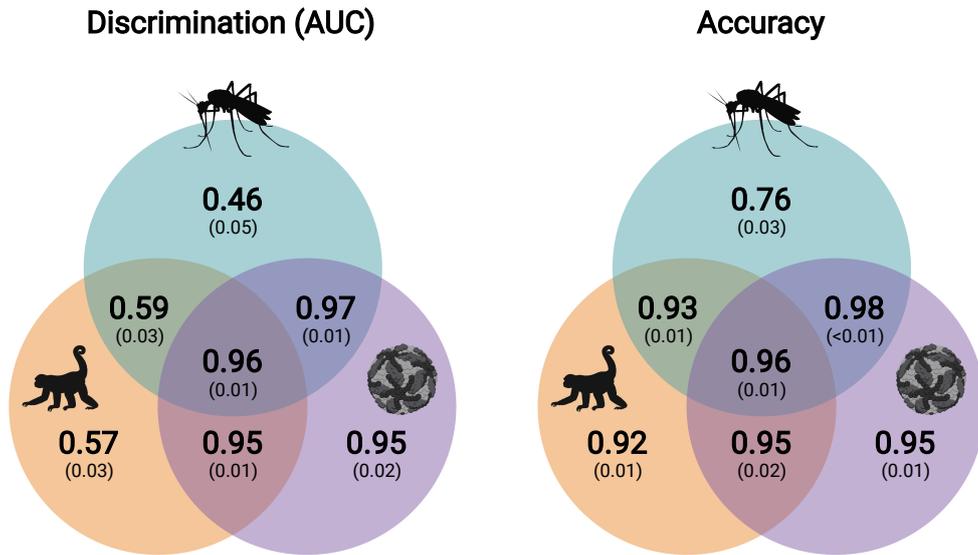


Figure 2: Model performance – quantified using AUC (left panel) and accuracy (right panel) – was highest when host, vector, and virus traits were included in the model (reported values are mean and standard deviation based on 20 model runs, assessing performance on a random 20% subset of the data). However, host-virus association model performance was not appreciably increased by the addition of 32 host trait covariates, suggesting that host-virus associations may be best predicted by considering information on the vector and the virus.

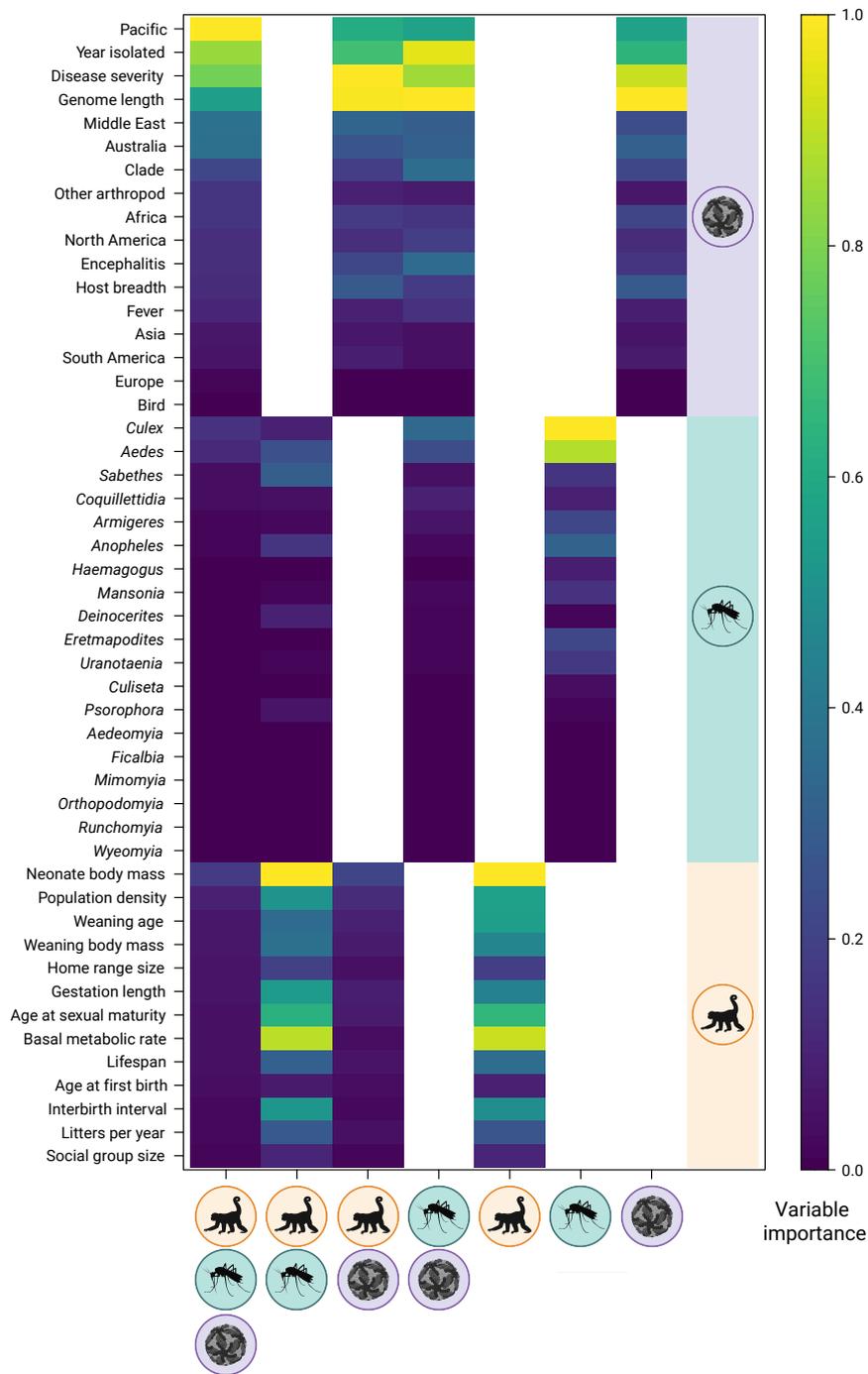


Figure 3: The relative importance of host (orange), vector (blue), and virus (purple) traits on predictive model performance. Each column corresponds to a different combination of these three trait groups, with the first column corresponding to the full model (as indicated at the bottom of each column using the glyphs). Variables are ordered based on the full model.

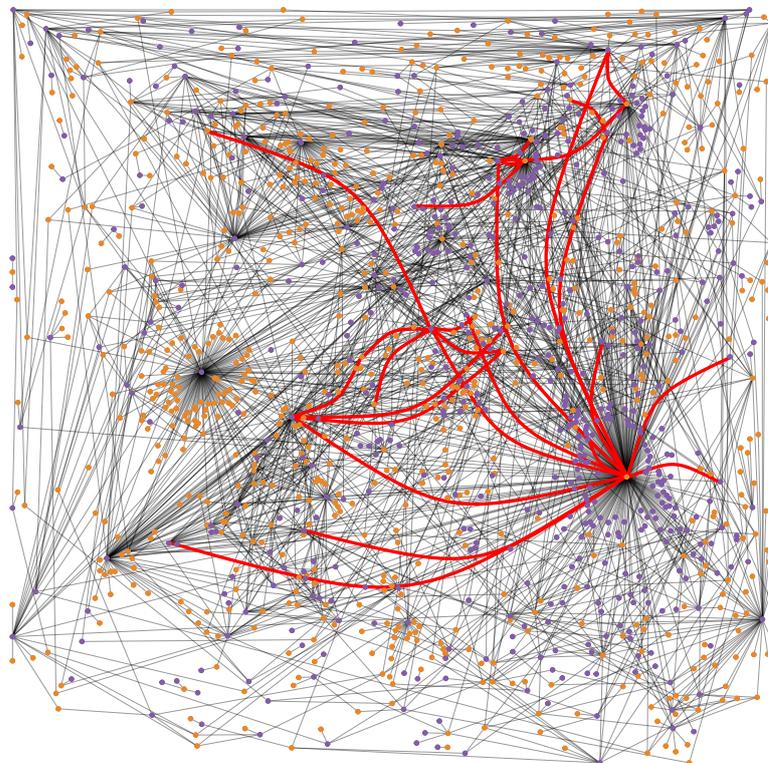


Figure 4: Full graph of host-virus associations (host species are in orange and viruses in purple), where links between host and virus species represent known associations. Red links are those which the full host-vector-virus cycle is known.

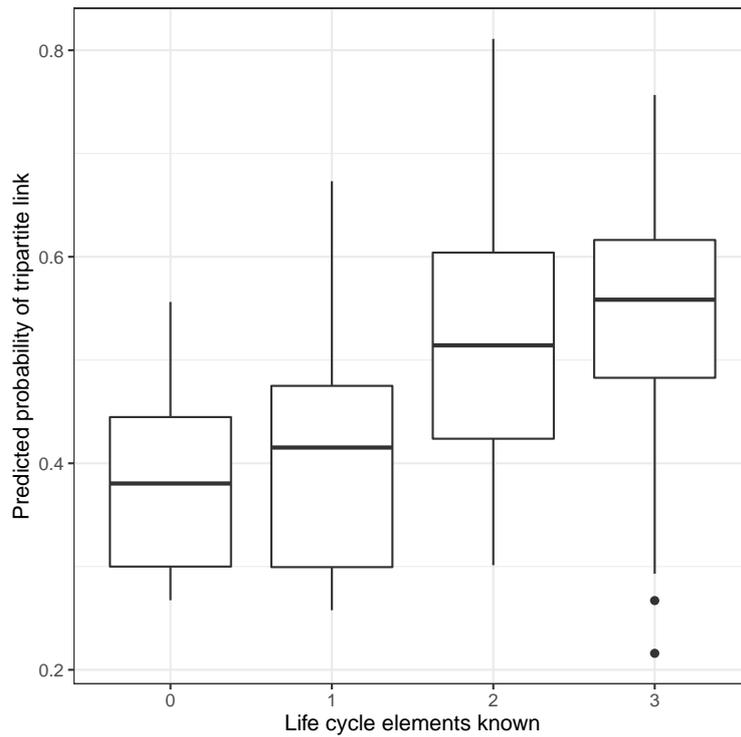


Figure 5: The tripartite model predicts a higher average probability for associations that have one or two links known (which are still not recorded as positive values in the training data) than those with no elements known to be possible. This suggests that the model is capable of more than just recapitulating the data, and is able to distinguish different levels of biological plausibility within unknown tripartite elements.

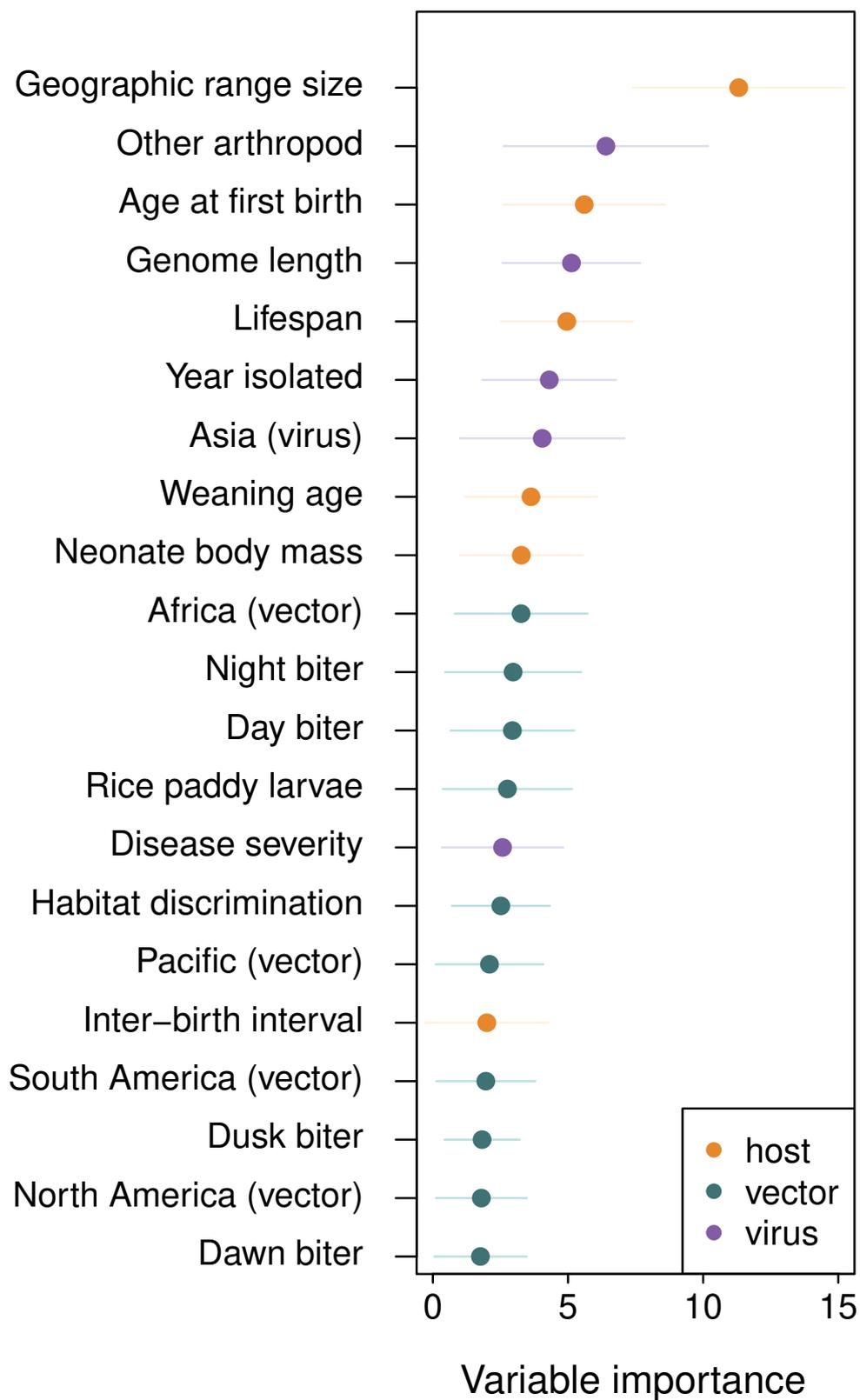


Figure 6: The relative importance of host (orange), vector (grey), and virus (blue) traits on predictive model performance in the tripartite model.

Tables

Table 1: **Top predicted epidemic cycles in humans.** All vectors are known to be human-biting; all viruses are known to be zoonotic based on either clinical or serological data.

<i>Host</i>	<i>Mosquito</i>	<i>Virus</i>	<i>Prob</i>
<i>H. sapiens</i>	<i>Culex pipiens</i>	Wesselsbron virus	0.81
<i>H. sapiens</i>	<i>Aedes aegypti</i>	West Nile virus	0.74
<i>H. sapiens</i>	<i>Aedes aegypti</i>	Japanese encephalitis virus	0.73
<i>H. sapiens</i>	<i>Culex pipiens</i>	Murray Valley encephalitis virus	0.73
<i>H. sapiens</i>	<i>Culex sitiens</i>	West Nile virus	0.72
<i>H. sapiens</i>	<i>Aedes scapularis</i>	West Nile virus	0.68
<i>H. sapiens</i>	<i>Mansonia uniformis</i>	West Nile virus	0.68
<i>H. sapiens</i>	<i>Anopheles coustani</i>	Wesselsbron virus	0.67
<i>H. sapiens</i>	<i>Culex pipiens</i>	Yellow fever virus	0.67
<i>H. sapiens</i>	<i>Aedes aegypti</i>	Ilheus virus	0.66
<i>H. sapiens</i>	<i>Aedes albopictus</i>	Ilheus virus	0.65
<i>H. sapiens</i>	<i>Anopheles hyrcanus</i>	Kokobera virus	0.62
<i>H. sapiens</i>	<i>Culex nigripalpus</i>	Wesselsbron virus	0.61
<i>H. sapiens</i>	<i>Aedes cantans</i>	Wesselsbron virus	0.61
<i>H. sapiens</i>	<i>Mansonia africana</i>	West Nile virus	0.61
<i>H. sapiens</i>	<i>Culex perexiguus</i>	Wesselsbron virus	0.60
<i>H. sapiens</i>	<i>Culex thalassius</i>	Wesselsbron virus	0.60
<i>H. sapiens</i>	<i>Culex gelidus</i>	West Nile virus	0.60
<i>H. sapiens</i>	<i>Culex annulirostris</i>	St. Louis encephalitis virus	0.59
<i>H. sapiens</i>	<i>Anopheles pharoensis</i>	West Nile virus	0.59

Table 2: **Top predicted enzootic cycles.** All mammals in the top 20 are either species found alongside humans (cows, sheep, pigs, and rats) or easily-sampled species from eastern North America (deer, raccoons, and possums).

<i>Host</i>	<i>Mosquito</i>	Virus	Prob
<i>Sus scrofa</i>	<i>Aedes albopictus</i>	West Nile virus	0.70
<i>Sus scrofa</i>	<i>Mansonia uniformis</i>	Wesselsbron virus	0.69
<i>Didelphis virginiana</i>	<i>Aedes aegypti</i>	Wesselsbron virus	0.67
<i>Didelphis virginiana</i>	<i>Aedes albopictus</i>	Wesselsbron virus	0.67
<i>Sus scrofa</i>	<i>Anopheles coustani</i>	Wesselsbron virus	0.66
<i>Sus scrofa</i>	<i>Culex quinquefasciatus</i>	Japanese encephalitis virus	0.65
<i>Procyon lotor</i>	<i>Culex tritaeniorhynchus</i>	West Nile virus	0.62
<i>Odocoileus virginianus</i>	<i>Anopheles pharoensis</i>	Wesselsbron virus	0.62
<i>Procyon lotor</i>	<i>Culex pipiens</i>	Japanese encephalitis virus	0.61
<i>Didelphis virginiana</i>	<i>Aedes aegypti</i>	West Nile virus	0.59
<i>Odocoileus virginianus</i>	<i>Aedes albopictus</i>	St. Louis encephalitis virus	0.56
<i>Bos taurus</i>	<i>Aedes vexans</i>	West Nile virus	0.56
<i>Bos taurus</i>	<i>Culex tritaeniorhynchus</i>	West Nile virus	0.55
<i>Procyon lotor</i>	<i>Culex annulirostris</i>	West Nile virus	0.55
<i>Macropus agilis</i>	<i>Culex tritaeniorhynchus</i>	Murray Valley encephalitis virus	0.54
<i>Bos taurus</i>	<i>Anopheles maculipennis</i>	Wesselsbron virus	0.54
<i>Ovis aries</i>	<i>Culex quinquefasciatus</i>	Ilheus virus	0.53
<i>Procyon lotor</i>	<i>Culex tarsalis</i>	West Nile virus	0.52
<i>Rattus rattus</i>	<i>Aedes aegypti</i>	Zika virus	0.51
<i>Odocoileus virginianus</i>	<i>Culex tritaeniorhynchus</i>	Banzi virus	0.51

515 Supplemental Material

516 Predicting the tripartite network of mosquito-borne disease

517 Trait data

518 Trait data were compiled from a variety of sources, with host trait data coming
 519 from PanTHERIA [21], and vector and virus trait data from Evans et al. 2017
 520 [13].

Table S1: Host, vector, and virus covariates considered in the models of host-virus (h-v column) and host-vector-virus (h-m-v column) associations. See the Pantheria documentation (<https://esapubs.org/archive/ecol/E090/184/metadata.htm>) for more information on host trait variables.

Taxa	Variable	Units	Definition	h-v	h-m-v
Host					
	Lifespan	days	Maximum observed lifespan	☑	☑
	Age at sexual maturity	days	Age at which individual is sexually mature	☑	☑
	Home range size	km ²	Area used by individual for daily tasks on average	☑	
	Gestation length	days	Period of time young are gestated	☑	☑
	Neonate body mass	grams	Average neonate body mass	☑	☑
	Population density	n / km ²	Number of individuals per unit area, on average	☑	
	Age at first birth	days	Age at which females give birth to their first litter	☑	☑
	Litters per year	n / year	Average number of litters per year	☑	
	Max lifespan	months	Longest observed lifespan	☑	☑
	Basal metabolic rate	mLO ₂ / hr	Individual metabolic rate	☑	

Interbirth interval	months	Period in between reproductive bouts	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Age at eye opening	days	Time when neonates open eyes	<input checked="" type="checkbox"/>	
Social group size	count	Number of individuals per social group	<input checked="" type="checkbox"/>	
Adult forearm length	mm	Length of adult forearm	<input checked="" type="checkbox"/>	
Dispersal age	days	Age at which young leave parents	<input checked="" type="checkbox"/>	
Neonate head-body length	mm	Body length of neonates	<input checked="" type="checkbox"/>	
Weaning age	days	Period of time when young stop weaning	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Weaning body mass	grams	Mass of young during weaning	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Vector				
Mosquito genus	numeric	Number of mosquito species of genus that vector a given virus		<input checked="" type="checkbox"/>
Human biter	1/0	Vector bites humans	<input checked="" type="checkbox"/>	
Host breadth	count	Number of host species bitten	<input checked="" type="checkbox"/>	
Non-primate mammals	1/0	Are non-primate mammals bitten	<input checked="" type="checkbox"/>	
Geographic range	count	Number of countries species collected	<input checked="" type="checkbox"/>	
Geographic location	-	Could include any or all of the following; Africa, Middle East, Australia, Pacific, Asia, Europe, North America, South America	<input checked="" type="checkbox"/>	
Biting behavior	-	Timing of biting behavior. Can be; dawn, day, dusk, and/or night	<input checked="" type="checkbox"/>	
Artificial container	1/0	Vector breeds in artificial containers	<input checked="" type="checkbox"/>	

Oviposition site	-	Larval site. Could include one or all of; treehole, container, pond, rockhole, marsh, swamp, ground pool, or rice paddy	<input checked="" type="checkbox"/>
Permanent habitat	1/0	Species uses permanent habitat	<input checked="" type="checkbox"/>
Habitat discrimination	count	number of habitat types	<input checked="" type="checkbox"/>
Urban preference	1/0	vector shows urban preference	<input checked="" type="checkbox"/>
Indoor preference	1/0	vector shows indoor preference	<input checked="" type="checkbox"/>
Viral range	count	Number species within genus to harbor virus	<input checked="" type="checkbox"/>

Virus

Average genome length	numeric	Length of viral genome	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Geographic location	-	Could include any or all of the following; Africa, Middle East, Australia, Pacific, Asia, Europe, North America, South America	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Clade	-	Viral clade (roman numerals)	<input checked="" type="checkbox"/>	
Year isolated	year	Virus isolation year	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Other arthropod	1/0	Vectored by other arthropods	<input checked="" type="checkbox"/>	
Host breadth	count	number of known hosts	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Encephalitis	1/0	Virus causes encephalitis	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Fever	1/0	Virus causes fever	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Disease severity	numeric	How severe is disease	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Bird host	1/0	Virus infects birds		<input checked="" type="checkbox"/>

What if we consider all vector-borne viruses?

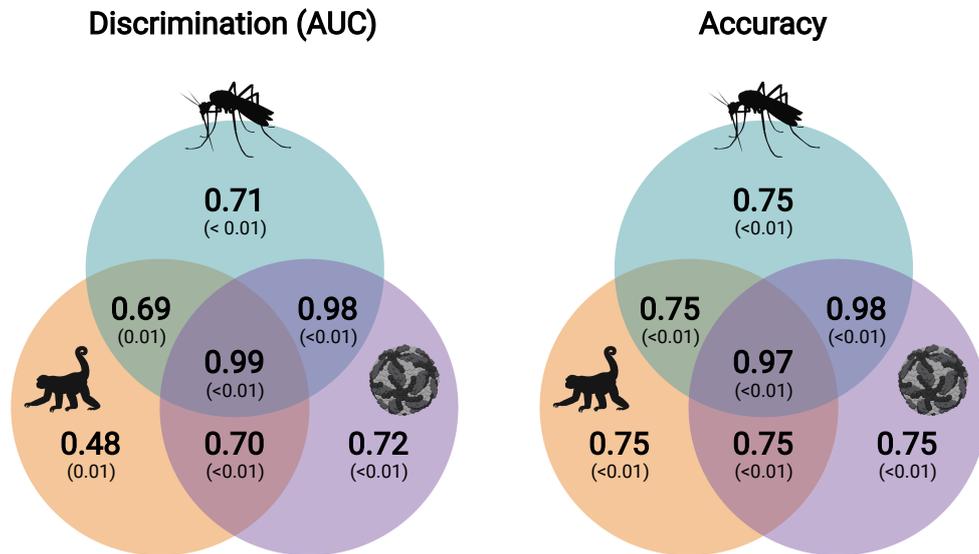


Figure S1: Model performance – quantified using AUC (left panel) and accuracy (right panel) – was highest when host, vector, and virus traits were included in the model (reported values are mean and standard deviation based on 20 model runs, assessing performance on a random 20% subset of the data). However, host-virus association model performance was not appreciably increased by the addition of vector data compared to just host and vector traits (AUC = 0.98).

522 **Different models and similar predictions**

523 When predicting host-virus associations, the different models had quite differ-
524 ent variable importance values, apart from obviously having different explanatory
525 variables. One question we had was whether models trained on different covari-
526 ates would not only have similar overall performance, but identify the same likely
527 host-virus associations as other models. To explore this graphically, we generated
528 a correlation matrix (Figure S2), where we find strong positive relationships be-
529 tween different model predictions. Interestingly, the least positive correlation was
530 from the full model, suggesting that the predictions from the full model differed
531 from models which consisted of nested subsets of the same features as in the full
532 model.

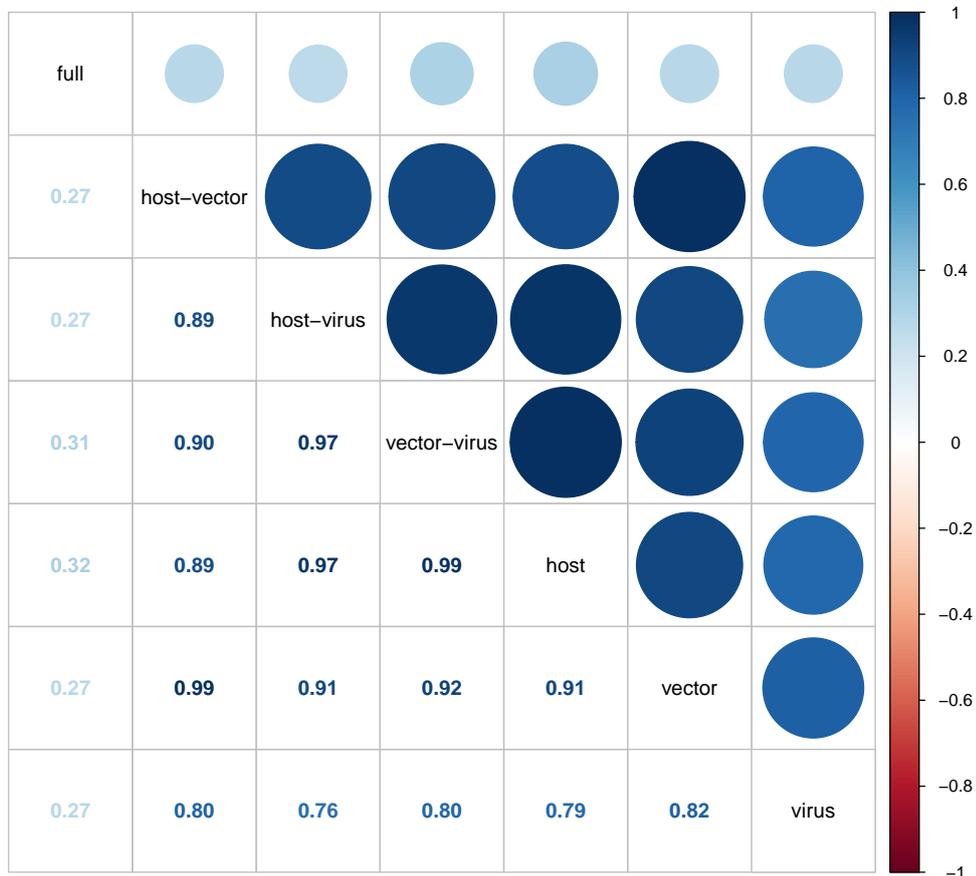


Figure S2: Correlation matrix between model predictions of the full set of subset models including different combinations of host, vector, and virus traits. The full model, including all traits, resulted in the predictions that were most weakly related to the other model predictions, though this model had similar performance as other models (see main text Figure 2). Lower triangle values and color scale correspond to Pearson's correlation coefficient values.