

1 **Ensemble Learning: Predicting Human Pathogenicity of Hematophagous**
2 **Arthropod Vector-Borne Viruses**

3

4 Huakai Hu^{1,2,#}, Chaoying Zhao^{3,1,#}, Meiling Jin⁴, Jiali Chen⁵, Xiong Liu¹, Hua Shi¹, Jinpeng
5 Guo¹, Changjun Wang^{3,1,*}, Yong Chen^{1,2,*}

6

7 ¹School of Public Health, China Medical University, Shenyang, Liaoning province, China

8 ²China Chinese PLA Center for Disease Control and Prevention, Beijing, China

9 ³School of Public Health, Zhengzhou University, Zhengzhou, Henan Province, China

10 ⁴Liaoning Provincial Center for Disease Control and Prevention, Shenyang, Liaoning,
11 People's Republic of China

12 ⁵School of Medicine, NanKai University, Tianjin, People's Republic of China.

13

14 Running title: Predicting Pathogenicity of Hematophagous Arthropod Vector-Borne Viruses

15

16 *These authors contributed equally to this work and are listed as co-first authors

17 Correspondence to Yong Chen chenyonger@126.com

18 Changjun Wang, science2008@hotmail.com

19

20 **Abstract**

21 Hematophagous arthropods occupy a pivotal role in ecosystems, serving as vectors for a wide
22 array of pathogens with significant implications for public health. Their capacity to harbor
23 and transmit viruses through biting actions creates a substantial risk of zoonotic spillover.

24 Despite the advancements in metagenomic approaches for virus discovery in vectors, the
25 isolation and cultivation of viruses still pose significant challenges, thereby limiting
26 comprehensive assessments of their pathogenicity. Here, we curated two datasets: one with
27 294 viruses, characterized by 37 epidemiological features, encompassing virus information
28 and host associations; the second with 71,622 sequences of hematophagous arthropod
29 vector-borne viruses, annotated with 33 sequence features. Two XGBoost models were

30 developed to predict arbovirus human pathogenicity—one integrating macroscopic
31 eco-epidemiological data, the other incorporating virus-related sequence features. The
32 macroscopic model identified non-vector host transmission as a key determinant, especially
33 involving Perissodactyla, Artiodactyla, and Carnivora Order. The sequence-based model
34 demonstrated that viral adhesion and viral invasion had distinct trends with consistent
35 increase and decrease in the likelihood of virus pathogenicity to humans, respectively. With
36 validated through an independent dataset, the model exhibited a congruous alignment with
37 documented pathogenicity outcomes. Together, the models offer a holistic framework for
38 assessing the pathogenic potential of viruses transmitted by hematophagous arthropods.

39

40

41 **Introduction**

42 Hematophagous arthropods, such as mosquitoes and ticks, play a pivotal role in ecosystems as
43 blood consumers and crucial disease vectors (Cuthbert et al., 2023; Touray et al., 2023). These
44 arthropods can harbor a myriad of pathogens, including bacteria, fungi, and viruses. Notably,
45 viral infections in these organisms are classified under the umbrella terms of Arthropod-Borne
46 Viruses (arboviruses) and insect-specific viruses (ISVs) (Calisher & Higgs, 2018; Gould et al.,
47 2017; Nouri et al., 2018; Zhao et al., 2022). The potential for these arthropods to harbor and
48 disseminate a diverse array of pathogens poses a grave threat to both human and animal
49 health, with the ominous potential to trigger outbreaks and result in a substantial number of
50 annual fatalities (Batson et al., 2021; Roth et al., 2018). Vector-borne diseases contribute
51 significantly to infectious diseases (Chala & Hamde, 2021), with notable arboviruses
52 including the Zika virus (Khongwichit et al., 2023; Weaver et al., 2018), Japanese encephalitis
53 virus (JEV) (Kampen & Werner, 2014), and the incessant menace of Dengue virus (DENV)
54 (Fournet et al., 2023).

55 In recent years, propelled by the widespread adoption of Viral Metagenomics sequencing
56 technologies, the identification of a wide range of established and emerging viruses within
57 hematophagous vectors, such as mosquitoes and ticks, has become feasible (Ni et al., 2023; X.
58 Yang et al., 2023). This technological progress presents an unprecedented opportunity to

59 comprehensively explore the distribution and transmission patterns of arboviruses and ISVs
60 across a spectrum of hosts, including both vectors and non-vectors. Such advancements are
61 crucial for supporting early warning systems, facilitating the anticipation and mitigation of
62 disease spread before its onset (Birnberg et al., 2020; Brinkmann et al., 2016). Despite these
63 achievements in viral metagenomics, current bioinformatic methods for virus recognition still
64 face limitations (Fang et al., 2019). Accurate identification of a significant number of
65 unknown contigs remains challenging. Even when identifying known or novel viruses, the
66 direct isolation and cultivation of these viruses from vectors proves to be formidable tasks,
67 hindering in-depth exploration of their pathogenesis and immune response (Lewis et al.,
68 2021).

69 In general, the close phylogenetic relatedness among viruses can offer insights into their
70 potential for human infectivity, as closely related viruses are generally presumed to share
71 common phenotypes and host ranges (Geoghegan & Holmes, 2018). However, despite being
72 a common rule of thumb for virus risk assessment, the extent to which evolutionary proximity
73 to viruses with known human infectivity accurately predicts zoonotic potential remains
74 unexamined in the current literature (Behl et al., 2022). Furthermore, the specific model is
75 designed to be trained on sequence features of closely related viruses (i.e., strains of the same
76 species) to discern viruses with human infectivity (Zhang et al., 2019). Unfortunately, this
77 method often overlooks critical functional characteristics of the viral genome, resulting in a
78 model that is less inclined to identify universally applicable pathogenic features across
79 diverse viruses. Consequently, predictions derived from such a model may be highly
80 susceptible to substantial biases (Mollentze et al., 2021).

81 The epidemiological characteristics of virus transmission encompass not only the virus itself
82 and information about its vector host (Zaid et al., 2021; Y.-J. S. Huang et al., 2019a; Viglietta
83 et al., 2021) but also factors such as geographical and climatic variations, as well as
84 interactions with non-vector hosts (Ciota & Keyel, 2019; Conway et al., 2014; Forrester et al.,
85 2014; Tabachnick, 2016). Moreover, for specific viruses, their nucleotide sequence
86 information may reflect actual pathogenic details (Bartoszewicz, Genske, et al., 2021).
87 Therefore, through a comprehensive analysis that integrates both macroscopic and
88 microscopic perspectives, our objective is to identify the epidemiological features and viral

89 sequence characteristics that have the greatest impact on the potential pathogenicity to
90 humans.

91 Based on the global data of arthropod-borne virus compiled by Huang et al. (Y. Huang et al.,
92 2023) as a foundation, we carefully curated the contents to extract pertinent information
93 concerning hematophagous arthropod-borne viruses. Additionally, to augment our analysis,
94 we utilized SeqScreen for insightful functional details of the viral sequences (Balaji et al.,
95 2022). Employing the XGBoost algorithm with ensemble learning, we developed both
96 regression and classification prediction models. This facilitated the identification of factors
97 with the most significant impact on human pathogenicity and enabled the construction of
98 ensemble learning for predicting the pathogenicity of virus sequences carried by
99 hematophagous arthropods.

100

101 **Materials and methods**

102 **Database restructuring and epidemiological feature retrieval**

103 The initial dataset comprised 101,094 virus sequences sourced from NCBI, spanning the
104 period from March 11, 1991, to January 28, 2023 (Y. Huang et al., 2023). To enhance the
105 reliability and specificity of our analysis, a stringent screening process was applied,
106 systematically excluding records lacking host information, sampling location details, and
107 those with ambiguous vector-host relationships. It is important to note that this dataset
108 excludes data from Antarctica. Subsequently, we identified 11 species of hematophagous
109 arthropods, including mosquitoes and ticks (Table supplement 1), while excluding
110 non-blood-feeding species such as *Tipulidae* and *Chironomidae*. Following this, we
111 systematically screened the entire dataset, retaining records exclusively related to hosts
112 classified as hematophagous arthropods. This refined dataset, derived through meticulous
113 curation, forms the foundation for our research, ensuring the integrity and accuracy of
114 subsequent analyses. Nevertheless, due to inaccuracies in the classification of vectors within
115 the database, a Python script was developed to scrape taxonomic directory. This script
116 retrieved detailed order, family, and genus information for each hematophagous arthropod and
117 non-vector host classification. To reveal the distinct composition of non-vector hosts, host

118 counts underwent logarithmic transformation (Figure 1A). For a more comprehensive
119 presentation, host classifications with fewer than 100 occurrences were amalgamated into an
120 “others” category, resulting in a total of 10 host classifications (Figure 1C).

121 In terms of additional epidemiological features, Köppen climate classification data for each
122 vector were acquired based on their discovery locations. This information was sourced from
123 both the Weather and Climate website (<https://weatherandclimate.com/>) and the Mindat
124 website (<https://www.mindat.org/>). Concurrently, continental data for each country were
125 obtained from the World Population Review (<https://worldpopulationreview.com/continents>),
126 and Baltimore classification data were sourced from the International Committee on
127 Taxonomy of Viruses (ICTV) (<https://ictv.global/report/genome>).

128 **The development of a regression model for macroscopic characteristics**

129 Firstly, among the 8,468 datasets in this study, We employed an R script to transform it into a
130 dataset comprising 294 distinct virus types, each characterized by 37 unique features (Table 1).
131 Subsequently, we utilized the XGBoost ensemble learning model to establish regression
132 models. The dataset was divided into training and validation sets at a ratio of 7.5:2.5. Given
133 the balanced ratio of positive to negative samples (1:1) in the model's database, addressing
134 imbalance was not deemed necessary. The training set was employed to train the model based
135 on the specified parameters (Table supplement 2). After determining the optimal number of
136 iterations through 10-fold cross-validation, we proceeded to construct the final model using
137 this identified count.

138 **Development and validation of a macroscopic features classification model**

139 In the microscopic pathogenicity classification model, we annotated the aforementioned
140 database uniformly using SeqScreen, resulting in a total of 71,622 virus sequences. After
141 excluding viruses from hosts submitted to NCBI after 2022, we obtained a final dataset of
142 71,593 sequences for this model. Due to the imbalance in positive and negative samples in the
143 database (positivity rate of 79.3%), we adjusted the sample sampling rate to balance the
144 dataset. The specific parameters employed in this model are meticulously detailed in Table
145 supplement 3. Utilizing the training set, the model was trained in accordance with these
146 parameters, determining the optimal iteration count through rigorous 5-fold cross-validation.
147 Subsequently, the final model was constructed utilizing the identified optimal iteration count.

148 To constitute an additional validation dataset, we retrieved Ebinur Lake Virus with arthropods
149 as hosts from NCBI, incorporating these samples with those previously excluded. Following a
150 consistent application of the specified parameters, we trained the model using functional
151 features from the entire dataset. Subsequently, predictions were generated on the additional
152 validation dataset. The obtained results underwent a meticulous comparative analysis with
153 findings from established pathogenic studies.

154

155

156 **Results**

157 **Global overview of hematophagous vector-virus distribution, diversity, and host** 158 **interactions**

159 This study has curated a comprehensive dataset of 8,468 hematophagous vector-virus pairs,
160 shedding light on their geographical distribution, diversity, and interactions with hosts. In
161 terms of distribution, these vectors were classified into two principal classes: *Insecta* and
162 *Arachnida*, spanning seven distinct families (Figure 1A). The records cover all six continents
163 except Antarctica, spanning across 102 countries globally and representing 24 diverse climate
164 types (Figure 1B). Regarding diversity, among the hematophagous vectors, *Culicidae* (64%,
165 5,445 sequences) predominates, constituting over half of all records, followed by *Ixodidae*
166 (32%, 2,703 sequences). Globally, the United States exhibits the highest diversity and
167 abundance of vectors, hosting five distinct families, followed by China with four. In terms of
168 virus records associated with vectors, the United States (1,977) leads the list, followed by
169 Russia, China, and Japan.

170 Turning to non-vector hosts, the dataset includes an additional 54,789 pairs of non-vector
171 hosts and viruses, with the non-vector hosts categorized into 15 groups. Among these hosts,
172 humans are the most prevalent, accounting for 40,078 records, followed by *Artiodactyla* and
173 *Aves*, constituting nearly 20% of the total (Figure 1C). The interactions between viruses and
174 non-vector hosts are distinct. The majority of viruses are associated with a single host.
175 Notably, West Nile virus (WNV) and Tick-borne encephalitis virus (TBEV), both belonging
176 to the *Flaviviridae*, exhibit the most widespread cross-host transmission, being detected in

177 nine non-vector host species. Moreover, as viruses expand their capacity to infect a wide
178 range of non-vector hosts, a noticeable reduction in viral diversity is observed. Specifically,
179 viruses capable of infecting only one host encompass 10 distinct virus families, while those
180 exhibiting infectivity across two to four hosts are confined to five families. Remarkably,
181 viruses with the ability to infect five, six, or seven hosts are prominently represented by
182 families such as *Flaviviridae* and *Togaviridae*. Among viruses capable of infecting seven
183 hosts, Dabie bandavirus stands out as a unique case. Belonging to the *Phenuiviridae*, this
184 virus is predominantly found in Asia (China, Japan, and South Korea). Infection with Dabie
185 bandavirus poses a significant health risk, causing a severe febrile illness accompanied by
186 thrombocytopenia, known as Severe Fever with Thrombocytopenia Syndrome (SFTS),
187 leading to its alternate nomenclature as the SFTS virus. The Japanese encephalitis virus
188 exhibits the highest degree of cross-vector host diversity, being detectable in three distinct
189 vector families: *Culicidae*, *Ixodidae*, and *Ceratopogonidae*. The previously mentioned WNV
190 and TBEV demonstrate transmission capabilities across both vector families, *Culicidae* and
191 *Ixodidae*.

192

193

194 **Pathogenicity of hematophagous arthropod vector-borne viruses: a macroscopic** 195 **regression analysis of epidemiological characteristics**

196 Through transforming the mentioned database and enhancing it with additional
197 epidemiological characteristics, we constructed a comprehensive dataset for the model. This
198 dataset comprises 294 distinct viruses, each characterized by 37 diverse features, broadly
199 categorized into viral characteristics, vector host features, and non-vector host features. To
200 unpack the crucial factors underlying human pathogenicity, we constructed and rigorously
201 trained an XGBoost model. This model leverages human infection status as the dependent
202 variable and incorporates 36 diverse features as independent variables, pinpointing the key
203 determinants of human infection. The model exhibits robust performance on the testing set,
204 with minimal prediction errors reflected in low MSE (0.01) and MAE (0.05) values,
205 highlighting its accuracy. Additionally, high R^2 (94.20%) and Explained Variance (94.29%)
206 values underscore the model's comprehensive ability to explain the variance in the dependent

207 variable.

208 The detection of viruses in non-vector hosts significantly influences human pathogenicity

209 (Figure 2), surpassing the impact of both vector hosts and the viral agents themselves.

210 Notably, the characteristics “Cross_host”, representing the total diversity of non-vector hosts

211 in which the virus has been detected, carries a weight of 52 in the model. This underscores the

212 critical role of the diversity of non-vector hosts in determining human pathogenicity.

213 Specifically, when considering potential human-pathogenicity, the order of importance is as

214 follows: *Perissodactyla*, *Artiodactyla Carnivora* and *Aves*. The higher the diversity of virus

215 detections in these animals, the greater the likelihood of the virus being pathogenic to humans.

216 After non-vector host factors, the subsequent important set of characteristics relates to the

217 vector hosts. Among these, “Cross_vector_g”, representing interspecies transmission among

218 diverse vector genus, emerges as the most critical factor. If a virus can propagate within

219 diverse vector genus, there is a substantial likelihood of viral spillover. The third set of

220 characteristics relates to the intrinsic characteristics of the virus itself. The closer the viral

221 phylogenetic relationship, the higher the likelihood of inducing similar immune responses,

222 thereby leading to diseases.

223

224 **Relationship between viral genomic function and human pathogenicity: a microscopic**

225 **machine learning approach**

226 In our research, we employed SeqScreen to functionally annotate all viral sequences in our

227 comprehensive dataset. After excluding sequences without successful annotations, our refined

228 dataset comprised 71,622 arboviruses and ISV sequences, each accompanied by their

229 respective functional, host, and pathogenic features. The largest category within our dataset

230 consists of mosquito-borne arboviruses, with Dengue virus 1 (9,194 sequences), Dengue virus

231 2 (8,999 sequences), and West Nile virus (4,656 sequences) being the most prevalent.

232 Tick-borne arboviruses, including African swine fever virus (3,915 sequences) and

233 Crimean-Congo hemorrhagic fever orthonairovirus (3,771 sequences) closely follow in

234 quantity.

235 Our functional annotation revealed a total of 10 distinct pathogenic features. The results

236 indicated that “viral adhesion” is the most prevalent function, accounting for 62% (44,482

237 sequences). This function facilitates virus adhesion to host cells, initiating infection and
238 paving the way for subsequent invasion and replication. Following closely are the “viral
239 counter-signaling” (49%) and “host xenophagy” (47%), which are typically associated with
240 immune evasion. These mechanisms enable the virus to survive, replicate within host cells,
241 and successfully transmit to other cells (Table 1).

242 Among known non-pathogenic viruses to humans, “viral invasion” stands out as the most
243 prevalent function, despite its relatively lower overall count compared to other functions.
244 Notably, within these viruses, the primary hosts targeted are hematophagous arthropod
245 vectors, with non-vector hosts predominantly represented by *Artiodactyla* and *Aves* (Figure
246 3A). Conversely, within the known human-pathogenic viruses, “viral adhesion” ranks as the
247 most prevalent function in terms of annotation quantities. In this context, excluding human
248 hosts, hematophagous arthropod vectors continue to be predominant, followed by *Aves*
249 (Figure 3B). This observation suggests a potential genomic similarity in the pathogenicity of
250 arboviruses among humans, hematophagous arthropod vectors and *Aves*.

251 We developed a binary classification XGBoost model using 33 features, which included
252 functional annotations for all viruses in the database and viral size (length of virus). The
253 model's dependent variable denotes whether a virus is pathogenic to humans. After excluding
254 the viruses in the extra validation dataset, the remaining viruses in the database were allocated
255 to a training set and a testing set in a 7.5:2.5 ratio. While achieving a high accuracy (95.36%),
256 we incorporated additional metrics, such as Precision (97.57%), Recall (96.55%), and F1
257 score (97.06%), for a more nuanced assessment. Furthermore, we generated an ROC curve
258 (Figure 4A) and a confusion matrix (Figure 4B) to gain a holistic view of the model's
259 strengths and weaknesses.

260 The model's results clearly demonstrate that, in terms of average gain, “viral adhesion”
261 exhibits the highest value, significantly enhancing the model's predictive accuracy. “Host
262 xenophagy” and “viral invasion” closely follow suit (Figure 5A). Regarding the model's
263 coverage, “viral invasion” and “host ubiquitin” occupy the top two positions due to their
264 capability to impact a wide array of viral sequences (Figure 5B). In terms of the model's
265 weights, the size of the viral sequence takes precedence over other features, indicating its
266 frequent utilization in the model construction and its vital role in making supplementary

267 assessments on the virus pathogenicity based on functional insights (Figure 5C).

268 In our conclusive analysis, we employed SHAP (SHapley Additive exPlanations) to gain
269 deeper insights into the individual feature contributions to the model. Notably, the top-ranking
270 feature— “viral size” —does not exhibit a discernible trend in pathogenicity to humans.
271 However, other features reveal intriguing patterns. Specifically, both “viral adhesion” and
272 “host xenophagy”, although slightly less significant than size, individually demonstrate
273 distinct trends: viral sequences annotated with either of these functions consistently increase
274 the likelihood of virus pathogenicity in humans. Conversely, “viral invasion” demonstrates an
275 inverse relationship, wherein sequences possessing this trait tend to reduce the probability of
276 virus pathogenicity. The majority of remaining features, on the other hand, positively
277 correlate with pathogenicity. In summary, most features contribute towards determining the
278 likelihood of virus pathogenicity in humans (Figure 6).

279 To delve deeper into the intricate interactions among these features, we conducted a thorough
280 analysis. Our results highlight that among all features, “viral counter signaling” exhibits the
281 most significant interaction with viral size. However, its impact on pathogenicity remains
282 inconclusive, lacking a definitive directional trend (Figure 7A). Additionally, we observed a
283 noteworthy interaction between “host xenophagy” and “viral adhesion”. The concurrent
284 presence of these features substantially enhances the virus's pathogenicity towards humans
285 (Figure 7B). Interestingly, “viral invasion” demonstrates a strong but contrasting interaction
286 with “viral counter signaling”. Specifically, “viral counter signaling” seems to function as a
287 protective factor against human pathogenicity when “viral invasion” is present, leading to a
288 reduced likelihood of the virus being pathogenic to humans (Figure 7C).

289 In our comprehensive analysis of interactions across all features, significant insights emerged.
290 While viral size lacks a clear discernible trend in its interaction with other features, the
291 interplay of 'host xenophagy' with both “viral adhesion” and 'viral counter signaling' guides
292 the model toward non-pathogenicity predictions, acting as a protective feature (Figure S1).

293 To assess real-world performance, we compiled an additional dataset consisting of 29 viruses
294 carried by hematophagous vectors, submitted after 2022. This dataset comprises 24 strains of
295 SFTS virus, 3 of Restan viruses, 3 of Tataguine viruses, 1 of Japanese encephalitis virus
296 (JEV), and 1 of Nairobi sheep disease virus (NSDV). Furthermore, 25 sequences of Ebinur

297 Lake virus (EBIV) borne by arthropods were downloaded from the NCBI, resulting in a total
298 of 54 virus sequences. To ensure independence, the content of the additional dataset was
299 excluded from the original model database. The train and test datasets were merged to train
300 the ultimate model. Validation was then conducted using the additional dataset. Model
301 predictions indicated that all sequences of SFTSV, a specific strain of JEV, a particular variant
302 of Tataguine virus, and one isolate of Ebinur Lake Virus potentially exhibit pathogenicity to
303 humans.

304

305

306 **Discussion**

307 The intricate interplay between hematophagous arthropods and the viruses they harbor forms
308 a dynamic ecosystem with profound implications for public health. Ticks and mosquitoes,
309 integral components of ecosystems, serve as potent disease vectors capable of transmitting
310 various pathogens, including arboviruses and ISVs. Recent advancements in Viral
311 Metagenomics sequencing technologies have significantly transformed our investigation of
312 the virome within hematophagous vectors, providing unparalleled access to comprehensive
313 viral genetic information. Although technological advancements have been made,
314 bioinformatic methods for virus recognition still face inherent limitations, particularly in
315 identifying unknown contigs, which hinders comprehensive virome characterization. The
316 isolation and cultivation of known or novel viruses from vectors remains challenging,
317 impeding in-depth exploration of their pathogenesis and immune response. The common
318 approach to assessing pathogenicity relies on viral phylogenetic analysis, assuming that
319 viruses with significant phylogenetic distance share similar pathogenic properties. However,
320 the extent to which phylogenetic relatedness accurately predicts the potential for zoonotic
321 diseases remains a critical aspect requiring deeper exploration. Existing models, often tailored
322 for closely related viruses, may inadvertently overshadow crucial functional characteristics,
323 introducing biases in predictions. To address these challenges, this study innovatively adopts
324 an ensemble learning algorithm in machine learning. Our objective is to comprehensively
325 explore the human-pathogenicity of viruses borne by hematophagous vectors, considering

326 both macro and micro-level characteristics, with the ultimate aim of identifying key viral
327 features associated with pathogenicity.

328 Based on the curated dataset, the distribution and abundance of hematophagous arthropods
329 suggest that the United States, Russia, and China harbor the highest number of vector insects,
330 notably mosquitoes and ticks, acting as primary carriers for arthropod-borne viruses (Y.-J. S.
331 Huang et al., 2019b; Wu et al., 2023). Among the viruses carried by these arthropods, RNA
332 viruses, specifically those classified as dsRNA under the third Baltimore group, dominate.
333 *Phlebotominae* and *Ceratopogonidae* share virus families, while distinct mosquito genera
334 harbor a diverse array of viruses, primarily belonging to the *Flaviviridae*, *Togaviridae*, and
335 *Peribunyaviridae* (Figure supplement 1). The abundance of viruses is influenced by various
336 factors, prompting a correlation analysis on the dataset. The results reveal strong correlations
337 for most viruses, excluding *Asfarviridae*, Zirqa virus, and Wallerfield virus, with six key
338 characteristics, including vector family and weather conditions. Within the community of
339 vector-borne viruses, *Flaviviridae*, *Togaviridae*, *Bunyaviridae*, *Rhabdoviridae*, and
340 *Phenuiviridae* are commonly co-detected and considered as core virome (Coatsworth et al.,
341 2022) (Figure supplement 2).

342 Within the cyclic dynamics of arboviruses, numerous factors intricately influence the
343 transmission to humans, impacting pathogenic outcomes. Variations in the intrinsic nature of
344 viruses yield diverse levels of human pathogenicity, commonly associated with phylogenetic
345 proximity. Moreover, an increased diversity of viruses within vectors may foster co-infection,
346 thereby facilitating viral evolution and spill-over (Vogels et al., 2019). Consistent with these
347 observations, our results highlight the significant influence of viral intrinsic factors on human
348 pathogenicity. These viruses primarily propagate diseases through hematophagous vectors,
349 predominantly mosquitoes, The species and behaviors of these vectors, collectively termed
350 “Vector capacity” (Conway et al., 2014), along with environmental shifts (Hermanns et al.,
351 2023; Weissenböck et al., 2010), play a pivotal role in shaping vector composition and
352 consequentially impact viral transmission. Notably, our research reveals that the impact of
353 vector hosts is equivalent to that of viral intrinsic factors. Both the diversity across vector
354 genera and the quantity of vector genera exert a substantial influence on human pathogenicity.
355 Interactions between viruses and non-vector hosts drive viral evolution (Sen et al., 2016),

356 with interspecies interactions being the primary driver for viral spill-over (Y.-J. S. Huang et
357 al., 2019a). The interplay between vectors and *Aves* hosts enables long-distance viral
358 transmission (Forrester et al., 2014), while interactions with vertebrates also emerge as pivotal
359 determinants (García-Romero et al., 2023; Golnar et al., 2014; Stephenson et al., 2019). Our
360 results align with these insights, highlighting that, beyond interspecies categories, viruses
361 infecting *Perissodactyla* and *Artiodactyla* pose the most significant risk for human
362 pathogenicity, increasing the likelihood of transmission to humans and subsequent disease
363 outbreaks.

364 Capitalizing on the ability to analyze viral genomic data for predicting viral pathogenicity to
365 humans (Bartoszewicz, Seidel, et al., 2021), we conducted an exhaustive investigation into
366 the genetic functionalities of arboviruses and ISVs within the database. Utilizing the ensemble
367 learning algorithm, we meticulously developed and trained a predictive model. Furthermore,
368 an additional dataset comprising arboviruses submitted to NCBI after 2022 was incorporated
369 for validation and prediction purposes. Our model demonstrates superior performance,
370 showcasing distinctive contributions from individual functional features that collectively
371 shape the overarching trend in viral pathogenicity. Specifically, “viral adhesion”, representing
372 a pivotal mechanism for viral infection and entry into host cells, emerges not only as the
373 predominant feature but also significantly enhances the overall performance of the model.
374 Empirical evidence affirms that the presence of this feature in viral sequences, subsequent to
375 transmission to humans by hematophagous vectors, consistently indicates an elevated risk of
376 pathogenicity. For instance, viruses within the *Flaviviridae*, such as DENV, WNV, and ZIKV
377 (Begum et al., 2019; Cruz-Oliveira et al., 2015; Faustino et al., 2019; Hasan et al., 2017;
378 Martins et al., 2019), utilize E and capsid proteins to enter receptor cells. Likewise, the
379 Chikungunya virus, a member of the *Togaviridae*, facilitates the fusion with receptor cells
380 through trimeric E1/E2 spikes (Ciota & Keyel, 2019).

381 While “viral invasion” plays a pivotal role in the initial phase of viral entry, it is notably
382 scarce among this dataset (Table 1). Furthermore, most viruses with this feature are presently
383 classified as non-pathogenic to humans (Figure 3). The high abundance of *Flavivirus* in the
384 dataset could explain the limited occurrence of “viral invasion”, given their unique infection
385 mechanisms that may not require this specific feature. Additionally, SeqScreen might not

386 have detected “viral invasion” in these sequences. Despite its relatively low occurrence in the
387 dataset, this feature exhibits the highest cover value in the XGBoost model, indicating its
388 significant impact on the model's performance. Interestingly, contrary to expectations, its
389 presence predominantly acts as a “protective factor” in predicting pathogenicity, as revealed
390 by SHAP explanations. This phenomenon may be attributed to the prevalence of this feature
391 among non-pathogenic viruses in our dataset. However, since we balanced the samples during
392 model training, it could also be a result of interactions among different features. A more
393 in-depth exploration of interactions related to this trait revealed a robust interplay with “viral
394 counter signaling”. When both features coexist, the model significantly leans towards
395 predicting non-pathogenicity in humans. Importantly, these two processes are not mutually
396 exclusive factors in actual virus infections. This observation implies potential distinctive
397 invasion mechanisms of arboviruses, indicating unconventional pathways for entering host
398 cells that facilitate immune evasion.

399 In this dataset, “viral counter signaling” and “host xenophagy” are prevalent features actively
400 enhancing virus pathogenicity and triggering host infection. They play a crucial role in the
401 pathogenicity to humans (Costa et al., 2013; King et al., 2020). Notably, “host xenophagy”,
402 similar to “viral adhesion”, significantly influences the model results. In terms of interactions,
403 it has the strongest interaction with “viral adhesion”, leading to a positive inclination towards
404 pathogenicity.

405 The feature of “size”, representing the length of viral sequences, while not directly associated
406 with pathogenic functions, plays a crucial role in refining the final results, as indicated by the
407 model's weight (Figure 5C). Training the model with only 33 functional features resulted in
408 unreliable accuracy (82%) and a high false-positive rate. However, the inclusion of “size”
409 substantially improved the model's performance. Notably, the influence of “size” on viral
410 pathogenicity lacks a discernible trend (Figure 6), resulting in predictive outcomes that tend
411 towards a more stochastic distribution. Interaction analysis revealed that “viral counter
412 signaling” has the strongest interaction with “size” (Figure 7). Even with these interactions,
413 determining the direction of pathogenicity remains challenging. In summary, 'size' appears to
414 fine-tune the model's final predictions. When combined with functional features, it facilitates
415 a more accurate assessment of the likelihood of pathogenicity to humans.

416 In the validation results using an additional dataset, we identified four viruses with the
417 potential to infect humans and induce diseases. Firstly, all sequence of Dabie bandavirus were
418 predicted to be pathogenic. These viruses collected from ticks in Miyazaki Prefecture, Japan,
419 exhibited a high degree of homology through phylogenetic analysis with a virus previously
420 isolated from an SFTS patient, providing strong evidence for its potential pathogenicity (Sato
421 et al., 2021). The model also predicted the pathogenicity of a strain of JEV, detected in
422 mosquitoes in the Qinghai-Tibet region of China (Li et al., 2011). Despite the elevated
423 altitude of the region, the presence of antibodies against the virus in both the indigenous
424 population and swine suggests a localized occurrence of virus transmission, thereby
425 challenging the initial presumption that the virus would not be prevalent at higher elevations.
426 Tataguine virus (Kapusinski et al., 2021), isolated from *Anopheles sp.* in Gambia, belongs to
427 the *Peribunyaviridae*. While its pathogenicity to humans remains inconclusive, there is a high
428 likelihood of infection symptoms if transmitted through hematophagous vectors. Among the
429 25 strains of Ebinur Lake Virus, one isolated from *Hyalomma marginatum* in the Volgograd
430 region of Russia in 2023 was predicted to be capable of infecting humans. This virus,
431 commonly found in China's prevalent vector host, *Culex modestus*, has been studied
432 extensively for its ability to infect BALB/c mice, resulting in pronounced clinical symptoms
433 (Zhao et al., 2020). Recent studies have substantiated the capacity of *Aedes aegypti* to serve as
434 a vector for such viruses (C. Yang et al., 2022). Although antibodies have been detected in
435 human serum samples, the lack of positive RT-PCR results prevents a conclusive
436 determination of the virus's ability to infect humans and induce diseases (Xia et al., 2020).
437 This aligns with the model's prediction, as the six viruses isolated from *Culex modestus*,
438 included in the model, are unlikely to be pathogenic to humans.

439 This study endeavors to leverage machine learning methodologies for discerning overarching
440 factors influencing the pathogenicity of hematophagous vector-borne viruses in humans. Our
441 developed predictive model, focused on gene function, has successfully demonstrated the
442 capability to predict virus pathogenicity in humans. However, it is crucial to acknowledge
443 certain limitations in our study. In the global dataset of vector-borne viruses, there exists an
444 uneven distribution, particularly with an overabundance of viruses such as DEV. This
445 imbalance may result in an unavoidable bias that impacts the accuracy of the model.

446 Furthermore, the selected machine learning algorithms, while effective, may not match the
447 efficacy of neural networks, posing challenges in optimizing for the current abundance of data.
448 Notably, variations in blood-feeding habits among hematophagous vectors were not
449 considered, which can significantly contribute to the spread of viruses. Different vector
450 species may exhibit distinct preferences and behaviors in their blood-feeding patterns,
451 influencing the transmission dynamics of viruses. Future research should incorporate these
452 behavioral nuances to provide a more comprehensive understanding of virus dissemination. In
453 summary, our model provides a novel perspective and serves as a valuable tool for the further
454 analysis of virus sequences, providing effective information for the monitoring and early
455 warning of hematophagous arthropod vector-borne transmission.

456 In this investigation, our primary objective is to discern both macroscopic and microscopic
457 factors influencing the risk of human pathogenicity in hematophagous vector-borne viruses.
458 Employing ensemble learning standpoint, we uncovered key characteristics associated with
459 viral pathogenicity from an epidemic perspective. Simultaneously, we delved into pivotal
460 functional features impacting human pathogenicity at a molecular level, with a specific focus
461 on the functional aspects of viral sequences. Moreover, we deploy our developed model to
462 forecast the human infectivity of viral sequences within an additional validation dataset. The
463 model's performance in predicting the pathogenicity of these viruses at the genetic level not
464 only enriches our comprehension of established and emerging virus risks but also broadens
465 the scope of hematophagous arthropods detection. Importantly, it contributes substantively to
466 the mitigation of present and future risks associated with vector-borne diseases.

467

468 **Acknowledgement**

469 The authors acknowledge the global open dataset shared by Huang et al and Xuan Li for
470 assistance with additional data collection. The laboratory is funded by a grant from National
471 Key Research and Development Program of China (2019YFC1200501). The funders had no
472 role in study design, data collection and interpretation, or the decision to submit the work for
473 publication.

474

475 **Data availability**

476 The data supporting the findings of this study are available upon reasonable request from the
477 author. Researchers interested in accessing the dataset for further exploration or verification
478 are encouraged to contact Huakai Hu at hhyu98@163.com for assistance. We are committed
479 to promoting transparency and collaboration in scientific research, and we welcome inquiries
480 regarding the data underlying our published results.

481

482 **Author contribution**

483 Huakai Hu, Idea Generation, Data Curation and Transformation, Model Development and
484 validation, Visualization, Writing - original; Chaoying Zhao, Conceptualization, Methodology,
485 Writing - original draft and review; Jiali Chen, Conceptualization, Methodology, Writing –
486 review and editing; Meiling Jin, Conceptualization, Writing – review and editing; Hua Shi,
487 Conceptualization, Writing – review and editing; Jinpeng Guo, Project administration,
488 Writing – review and editing; Changjun Wang, Conceptualization, Methodology, Writing –
489 review and editing; Yong Chen, Supervision, Funding acquisition, Project administration,
490 Writing – review and editing;

491

492

493 **References**

- 494 Balaji, A., Kille, B., Kappell, A. D., Godbold, G. D., Diep, M., Elworth, R. A. L., Q
495 ian, Z., Albin, D., Nasko, D. J., Shah, N., Pop, M., Segarra, S., Ternus, K. L., &
496 Treangen, T. J. (2022). SeqScreen: Accurate and sensitive functional screening of pa
497 thogenic sequences via ensemble learning. *Genome Biology*, 23(1), 133. [https://doi.or](https://doi.org/10.1186/s13059-022-02695-x)
498 [g/10.1186/s13059-022-02695-x](https://doi.org/10.1186/s13059-022-02695-x)
- 499 Bartoszewicz, J. M., Genske, U., & Renard, B. Y. (2021). Deep learning-based real-ti
500 me detection of novel pathogens during sequencing. *Briefings in Bioinformatics*, 22
501 (6), bbab269. <https://doi.org/10.1093/bib/bbab269>
- 502 Bartoszewicz, J. M., Seidel, A., & Renard, B. Y. (2021). Interpretable detection of no
503 vel human viruses from genome sequencing data. *NAR Genomics and Bioinformatics*,
504 3(1), lqab004. <https://doi.org/10.1093/nargab/lqab004>
- 505 Batson, J., Dudas, G., Haas-Stapleton, E., Kistler, A. L., Li, L. M., Logan, P., Ratnasi
506 ri, K., & Retallack, H. (2021). Single mosquito metatranscriptomics identifies vector
507 s, emerging pathogens and reservoirs in one assay. *eLife*, 10, e68353. [https://doi.org/](https://doi.org/10.7554/eLife.68353)
508 [10.7554/eLife.68353](https://doi.org/10.7554/eLife.68353)
- 509 Begum, F., Das, S., Mukherjee, D., & Ray, U. (2019). Hijacking the Host Immune Ce
510 lls by Dengue Virus: Molecular Interplay of Receptors and Dengue Virus Envelope.
511 *Microorganisms*, 7(9), Article 9. <https://doi.org/10.3390/microorganisms7090323>
- 512 Behl, A., Nair, A., Mohagaonkar, S., Yadav, P., Gambhir, K., Tyagi, N., Sharma, R.
513 K., Butola, B. S., & Sharma, N. (2022). Threat, challenges, and preparedness for fu
514 ture pandemics: A descriptive review of phylogenetic analysis based predictions. *Infe*
515 *ction, Genetics and Evolution: Journal of Molecular Epidemiology and Evolutionary*
516 *Genetics in Infectious Diseases*, 98, 105217. <https://doi.org/10.1016/j.meegid.2022.10>
517 [5217](https://doi.org/10.1016/j.meegid.2022.105217)
- 518 Birnberg, L., Temmam, S., Aranda, C., Correa-Fiz, F., Talavera, S., Bigot, T., Eloit,
519 M., & Busquets, N. (2020). Viromics on Honey-Baited FTA Cards as a New Tool f

- 520 or the Detection of Circulating Viruses in Mosquitoes. *Viruses*, 12(3), 274. <https://doi.org/10.3390/v12030274>
- 521
- 522 Brinkmann, A., Nitsche, A., & Kohl, C. (2016). Viral Metagenomics on Blood-Feeding
523 Arthropods as a Tool for Human Disease Surveillance. *International Journal of Mo*
524 *lecular Sciences*, 17(10), Article 10. <https://doi.org/10.3390/ijms17101743>
- 525 Calisher, C. H., & Higgs, S. (2018). The Discovery of Arthropod-Specific Viruses in
526 Hematophagous Arthropods: An Open Door to Understanding the Mechanisms of Ar
527 bovirus and Arthropod Evolution? *Annual Review of Entomology*, 63(1), 87–103. <https://doi.org/10.1146/annurev-ento-020117-043033>
- 528
- 529 Chala, B., & Hamde, F. (2021). Emerging and Re-emerging Vector-Borne Infectious Di
530 seases and the Challenges for Control: A Review. *Frontiers in Public Health*, 9. <https://www.frontiersin.org/articles/10.3389/fpubh.2021.715759>
- 531
- 532 Ciota, A. T., & Keyel, A. C. (2019). The Role of Temperature in Transmission of Zo
533 onotic Arboviruses. *Viruses*, 11(11), Article 11. <https://doi.org/10.3390/v11111013>
- 534 Coatsworth, H., Bozic, J., Carrillo, J., Buckner, E. A., Rivers, A. R., Dinglasan, R. R.,
535 & Mathias, D. K. (2022). Intrinsic variation in the vertically transmitted core viro
536 me of the mosquito *Aedes aegypti*. *Molecular Ecology*, 31(9), 2545–2561. <https://doi.org/10.1111/mec.16412>
- 537
- 538 Conway, M. J., Colpitts, T. M., & Fikrig, E. (2014). Role of the Vector in Arbovirus
539 Transmission. *Annual Review of Virology*, 1(1), 71–88. <https://doi.org/10.1146/annurev-virology-031413-085513>
- 540
- 541 Costa, V. V., Fagundes, C. T., Souza, D. G., & Teixeira, M. M. (2013). Inflammatory
542 and Innate Immune Responses in Dengue Infection: Protection versus Disease Indu
543 ction. *The American Journal of Pathology*, 182(6), 1950–1961. <https://doi.org/10.1016/j.ajpath.2013.02.027>
- 544
- 545 Cruz-Oliveira, C., Freire, J. M., Conceição, T. M., Higa, L. M., Castanho, M. A. R.

546 B., & Da Poian, A. T. (2015). Receptors and routes of dengue virus entry into the
547 host cells. *FEMS Microbiology Reviews*, 39(2), 155–170. <https://doi.org/10.1093/femsr>
548 e/fuu004

549 Cuthbert, R. N., Darriet, F., Chabrierie, O., Lenoir, J., Courchamp, F., Claeys, C., Rob
550 ert, V., Jourdain, F., Ulmer, R., Diagne, C., Ayala, D., Simard, F., Morand, S., & R
551 enault, D. (2023). Invasive hematophagous arthropods and associated diseases in a c
552 hanging world. *Parasites & Vectors*, 16(1), 291. <https://doi.org/10.1186/s13071-023-05>
553 887-x

554 Fang, Z., Tan, J., Wu, S., Li, M., Xu, C., Xie, Z., & Zhu, H. (2019). PPR-Meta: A t
555 ool for identifying phages and plasmids from metagenomic fragments using deep lea
556 rning. *GigaScience*, 8(6), giz066. <https://doi.org/10.1093/gigascience/giz066>

557 Faustino, A. F., Martins, A. S., Karguth, N., Artilheiro, V., Enguita, F. J., Ricardo, J.
558 C., Santos, N. C., & Martins, I. C. (2019). Structural and Functional Properties of t
559 he Capsid Protein of Dengue and Related Flavivirus. *International Journal of Molec
560 ular Sciences*, 20(16), Article 16. <https://doi.org/10.3390/ijms20163870>

561 Forrester, N. L., Coffey, L. L., & Weaver, S. C. (2014). Arboviral Bottlenecks and Ch
562 allenges to Maintaining Diversity and Fitness during Mosquito Transmission. *Viruses*,
563 6(10), Article 10. <https://doi.org/10.3390/v6103991>

564 Fournet, N., Voiry, N., Rozenberg, J., Bassi, C., Cassonnet, C., Karch, A., Durand, G.,
565 Grard, G., Modenesi, G., Lakoussan, S.-B., Tayliam, N., Zatta, M., Gallien, S., inv
566 estigation team, Noël, H., Brichler, S., & Tarantola, A. (2023). A cluster of autocht
567 honous dengue transmission in the Paris region—Detection, epidemiology and contro
568 l measures, France, October 2023. *Euro Surveillace: Bulletin Europeen Sur Les Ma
569 ladies Transmissibles = European Communicable Disease Bulletin*, 28(49). [https://doi.](https://doi.org/10.2807/1560-7917.ES.2023.28.49.2300641)
570 [org/10.2807/1560-7917.ES.2023.28.49.2300641](https://doi.org/10.2807/1560-7917.ES.2023.28.49.2300641)

571 García-Romero, C., Carrillo Bilbao, G. A., Navarro, J.-C., Martin-Solano, S., & Saeger
572 man, C. (2023). Arboviruses in Mammals in the Neotropics: A Systematic Review t

- 573 o Strengthen Epidemiological Monitoring Strategies and Conservation Medicine. *Viru*
574 *ses*, 15(2), Article 2. <https://doi.org/10.3390/v15020417>
- 575 Geoghegan, J. L., & Holmes, E. C. (2018). The phylogenomics of evolving virus viru
576 lence. *Nature Reviews Genetics*, 19(12), Article 12. <https://doi.org/10.1038/s41576-018>
577 -0055-5
- 578 Golnar, A. J., Turell, M. J., LaBeaud, A. D., Kading, R. C., & Hamer, G. L. (2014).
579 Predicting the Mosquito Species and Vertebrate Species Involved in the Theoretical
580 Transmission of Rift Valley Fever Virus in the United States. *PLOS Neglected Tro*
581 *pical Diseases*, 8(9), e3163. <https://doi.org/10.1371/journal.pntd.0003163>
- 582 Gould, E., Pettersson, J., Higgs, S., Charrel, R., & de Lamballerie, X. (2017). Emergi
583 ng arboviruses: Why today? *One Health*, 4, 1–13. <https://doi.org/10.1016/j.onehlt.201>
584 7.06.001
- 585 Hasan, S. S., Miller, A., Sapparapu, G., Fernandez, E., Klose, T., Long, F., Fokine, A.,
586 Porta, J. C., Jiang, W., Diamond, M. S., Crowe Jr., J. E., Kuhn, R. J., & Rossmann,
587 M. G. (2017). A human antibody against Zika virus crosslinks the E protein to
588 prevent infection. *Nature Communications*, 8(1), Article 1. <https://doi.org/10.1038/ncomms14722>
589
- 590 Hermanns, K., Marklewitz, M., Zirkel, F., Kopp, A., Kramer-Schadt, S., & Junglen, S.
591 (2023). Mosquito community composition shapes virus prevalence patterns along an
592 thropogenic disturbance gradients. *eLife*, 12, e66550. <https://doi.org/10.7554/eLife.665>
593 50
- 594 Huang, Y., Wang, S., Liu, H., Atoni, E., Wang, F., Chen, W., Li, Z., Rodriguez, S., Y
595 uan, Z., Ming, Z., & Xia, H. (2023). A global dataset of sequence, diversity and bi
596 osafety recommendation of arbovirus and arthropod-specific virus. *Scientific Data*, 10
597 (1), 305. <https://doi.org/10.1038/s41597-023-02226-8>
- 598 Huang, Y.-J. S., Higgs, S., & Vanlandingham, D. L. (2019a). Arbovirus-Mosquito Vect

- 599 or-Host Interactions and the Impact on Transmission and Disease Pathogenesis of Ar
600 boviruses. *Frontiers in Microbiology*, *10*. [https://www.frontiersin.org/articles/10.3389/f](https://www.frontiersin.org/articles/10.3389/fmicb.2019.00022)
601 [micb.2019.00022](https://www.frontiersin.org/articles/10.3389/fmicb.2019.00022)
- 602 Huang, Y.-J. S., Higgs, S., & Vanlandingham, D. L. (2019b). Emergence and re-emerg
603 ence of mosquito-borne arboviruses. *Current Opinion in Virology*, *34*, 104–109. <https://doi.org/10.1016/j.coviro.2019.01.001>
604 [//doi.org/10.1016/j.coviro.2019.01.001](https://doi.org/10.1016/j.coviro.2019.01.001)
- 605 Kampen, H., & Werner, D. (2014). Out of the bush: The Asian bush mosquito *Aedes*
606 *japonicus japonicus* (Theobald, 1901) (Diptera, Culicidae) becomes invasive. *Parasit*
607 *es & Vectors*, *7*, 59. <https://doi.org/10.1186/1756-3305-7-59>
- 608 Kapuscinski, M. L., Bergren, N. A., Russell, B. J., Lee, J. S., Borland, E. M., Hartm
609 an, D. A., King, D. C., Hughes, H. R., Burkhalter, K. L., Kading, R. C., & Stengl
610 ein, M. D. (2021). Genomic characterization of 99 viruses from the bunyavirus fami
611 lies Nairoviridae, Peribunyaviridae, and Phenuiviridae, including 35 previously unseq
612 uenced viruses. *PLoS Pathogens*, *17*(3), e1009315. [https://doi.org/10.1371/journal.ppat.](https://doi.org/10.1371/journal.ppat.1009315)
613 [1009315](https://doi.org/10.1371/journal.ppat.1009315)
- 614 Khongwichit, S., Chuchaona, W., Vongpunsawad, S., & Poovorawan, Y. (2023). Molec
615 ular epidemiology, clinical analysis, and genetic characterization of Zika virus infecti
616 ons in Thailand (2020-2023). *Scientific Reports*, *13*(1), 21030. [https://doi.org/10.1038/](https://doi.org/10.1038/s41598-023-48508-4)
617 [s41598-023-48508-4](https://doi.org/10.1038/s41598-023-48508-4)
- 618 King, C. A., Wegman, A. D., & Endy, T. P. (2020). Mobilization and Activation of t
619 he Innate Immune Response to Dengue Virus. *Frontiers in Cellular and Infection M*
620 *icrobiology*, *10*. <https://www.frontiersin.org/articles/10.3389/fcimb.2020.574417>
- 621 Lewis, W. H., Tahon, G., Geesink, P., Sousa, D. Z., & Ettema, T. J. G. (2021). Innov
622 ations to culturing the uncultured microbial majority. *Nature Reviews. Microbiology*,
623 *19*(4), 225–240. <https://doi.org/10.1038/s41579-020-00458-8>
- 624 Li, Y.-X., Li, M.-H., Fu, S.-H., Chen, W.-X., Liu, Q.-Y., Zhang, H.-L., Da, W., Hu, S.

- 625 -L., La Mu, S. D., Bai, J., Yin, Z.-D., Jiang, H.-Y., Guo, Y.-H., Ji, D. Z. D., Xu,
626 H.-M., Li, G., Mu, G. G. C., Luo, H.-M., Wang, J.-L., ... Liang, G.-D. (2011). Jap
627 anese Encephalitis, Tibet, China. *Emerging Infectious Diseases*, *17*(5), 934–936. <https://doi.org/10.3201/eid1705.101417>
628
- 629 Martins, A. S., Carvalho, F. A., Faustino, A. F., Martins, I. C., & Santos, N. C. (201
630 9). West Nile Virus Capsid Protein Interacts With Biologically Relevant Host Lipid
631 Systems. *Frontiers in Cellular and Infection Microbiology*, *9*. [https://www.frontiersin.
632 org/articles/10.3389/fcimb.2019.00008](https://www.frontiersin.org/articles/10.3389/fcimb.2019.00008)
- 633 Mollentze, N., Babayan, S. A., & Streicker, D. G. (2021). Identifying and prioritizing
634 potential human-infecting viruses from their genome sequences. *PLOS Biology*, *19*(9),
635 e3001390. <https://doi.org/10.1371/journal.pbio.3001390>
- 636 Ni, X.-B., Cui, X.-M., Liu, J.-Y., Ye, R.-Z., Wu, Y.-Q., Jiang, J.-F., Sun, Y., Wang,
637 Q., Shum, M. H.-H., Chang, Q.-C., Zhao, L., Han, X.-H., Ma, K., Shen, S.-J., Zha
638 ng, M.-Z., Guo, W.-B., Zhu, J.-G., Zhan, L., Li, L.-J., ... Cao, W.-C. (2023). Meta
639 virome of 31 tick species provides a compendium of 1,801 RNA virus genomes. *N
640 ature Microbiology*, *8*(1), 162–173. <https://doi.org/10.1038/s41564-022-01275-w>
- 641 Nouri, S., Matsumura, E. E., Kuo, Y.-W., & Falk, B. W. (2018). Insect-specific viruse
642 s: From discovery to potential translational applications. *Current Opinion in Virology*,
643 *33*, 33–41. <https://doi.org/10.1016/j.coviro.2018.07.006>
- 644 Roth, G. A., Abate, D., Abate, K. H., Abay, S. M., Abbafati, C., Abbasi, N., Abbasta
645 bar, H., Abd-Allah, F., Abdela, J., Abdelalim, A., Abdollahpour, I., Abdulkader, R.
646 S., Abebe, H. T., Abebe, M., Abebe, Z., Abejie, A. N., Abera, S. F., Abil, O. Z.,
647 Abraha, H. N., ... Murray, C. J. L. (2018). Global, regional, and national age-sex-s
648 pecific mortality for 282 causes of death in 195 countries and territories, 1980–2017:
649 A systematic analysis for the Global Burden of Disease Study 2017. *The Lancet*, *3
650 92*(10159), 1736–1788. [https://doi.org/10.1016/S0140-6736\(18\)32203-7](https://doi.org/10.1016/S0140-6736(18)32203-7)
- 651 Sato, Y., Mekata, H., Sudaryatma, P. E., Kirino, Y., Yamamoto, S., Ando, S., Sugimot

- 652 o, T., & Okabayashi, T. (2021). Isolation of Severe Fever with Thrombocytopenia S
653 yndrome Virus from Various Tick Species in Area with Human Severe Fever with
654 Thrombocytopenia Syndrome Cases. *Vector-Borne and Zoonotic Diseases*, *21*(5), 378
655 –384. <https://doi.org/10.1089/vbz.2020.2720>
- 656 Sen, R., Nayak, L., & De, R. K. (2016). A review on host–pathogen interactions: Cla
657 ssification and prediction. *European Journal of Clinical Microbiology & Infectious D*
658 *iseases*, *35*(10), 1581–1599. <https://doi.org/10.1007/s10096-016-2716-7>
- 659 Stephenson, E. B., Murphy, A. K., Jansen, C. C., Peel, A. J., & McCallum, H. (201
660 9). Interpreting mosquito feeding patterns in Australia through an ecological lens: A
661 n analysis of blood meal studies. *Parasites & Vectors*, *12*(1), 156. [https://doi.org/10.](https://doi.org/10.1186/s13071-019-3405-z)
662 [1186/s13071-019-3405-z](https://doi.org/10.1186/s13071-019-3405-z)
- 663 Tabachnick, W. J. (2016). Climate Change and the Arboviruses: Lessons from the Evo
664 lution of the Dengue and Yellow Fever Viruses. *Annual Review of Virology*, *3*(1), 1
665 25–145. <https://doi.org/10.1146/annurev-virology-110615-035630>
- 666 Touray, M., Bakirci, S., Ulug, D., Gulsen, S. H., Cimen, H., Yavasoglu, S. I., Simsek,
667 F. M., Ertabaklar, H., Ozbel, Y., & Hazir, S. (2023). Arthropod vectors of disease
668 agents: Their role in public and veterinary health in Turkiye and their control meas
669 ures. *Acta Tropica*, *243*, 106893. <https://doi.org/10.1016/j.actatropica.2023.106893>
- 670 Viglietta, M., Bellone, R., Blisnick, A. A., & Failloux, A.-B. (2021). Vector Specificit
671 y of Arbovirus Transmission. *Frontiers in Microbiology*, *12*. [https://www.frontiersin.o](https://www.frontiersin.org/articles/10.3389/fmicb.2021.773211)
672 [rg/articles/10.3389/fmicb.2021.773211](https://www.frontiersin.org/articles/10.3389/fmicb.2021.773211)
- 673 Vogels, C. B. F., Rückert, C., Cavany, S. M., Perkins, T. A., Ebel, G. D., & Grubaug
674 h, N. D. (2019). Arbovirus coinfection and co-transmission: A neglected public healt
675 h concern? *PLOS Biology*, *17*(1), e3000130. [https://doi.org/10.1371/journal.pbio.30001](https://doi.org/10.1371/journal.pbio.3000130)
676 [30](https://doi.org/10.1371/journal.pbio.3000130)
- 677 Weaver, S. C., Charlier, C., Vasilakis, N., & Lecuit, M. (2018). Zika, Chikungunya, an

- 678 d Other Emerging Vector-Borne Viral Diseases. *Annual Review of Medicine*, 69, 395
679 –408. <https://doi.org/10.1146/annurev-med-050715-105122>
- 680 Weissenböck, H., Hubálek, Z., Bakonyi, T., & Nowotny, N. (2010). Zoonotic mosquito
681 -borne flaviviruses: Worldwide presence of agents with proven pathogenicity and pot
682 ential candidates of future emerging diseases. *Veterinary Microbiology*, 140(3–4), 271
683 –280. <https://doi.org/10.1016/j.vetmic.2009.08.025>
- 684 Wu, Z., Zhang, M., Zhang, Y., Lu, K., Zhu, W., Feng, S., Qi, J., & Niu, G. (2023).
685 Jingmen tick virus: An emerging arbovirus with a global threat. *mSphere*, 8(5), e002
686 81-23. <https://doi.org/10.1128/msphere.00281-23>
- 687 Xia, H., Liu, R., Zhao, L., Sun, X., Zheng, Z., Atoni, E., Hu, X., Zhang, B., Zhang,
688 G., & Yuan, Z. (2020). Characterization of Ebinur Lake Virus and Its Human Serop
689 revalence at the China–Kazakhstan Border. *Frontiers in Microbiology*, 10. [https://ww
690 w.frontiersin.org/articles/10.3389/fmicb.2019.03111](https://www.frontiersin.org/articles/10.3389/fmicb.2019.03111)
- 691 Yang, C., Wang, F., Huang, D., Ma, H., Zhao, L., Zhang, G., Li, H., Han, Q., Bente,
692 D., Salazar, F. V., Yuan, Z., & Xia, H. (2022). Vector competence and immune res
693 ponse of *Aedes aegypti* for Ebinur Lake virus, a newly classified mosquito-borne or
694 thobunyavirus. *PLoS Neglected Tropical Diseases*, 16(7), e0010642. [https://doi.org/10.
695 1371/journal.pntd.0010642](https://doi.org/10.1371/journal.pntd.0010642)
- 696 Yang, X., Qin, S., Liu, X., Zhang, N., Chen, J., Jin, M., Liu, F., Wang, Y., Guo, J.,
697 Shi, H., Wang, C., & Chen, Y. (2023). Meta-Viromic Sequencing Reveals Virome C
698 haracteristics of Mosquitoes and *Culicoides* on Zhoushan Island, China. *Microbiology
699 Spectrum*, e02688-22. <https://doi.org/10.1128/spectrum.02688-22>
- 700 Zaid, A., Burt, F. J., Liu, X., Poo, Y. S., Zandi, K., Suhrbier, A., Weaver, S. C., Tex
701 eira, M. M., & Mahalingam, S. (2021). Arthritogenic alphaviruses: Epidemiological
702 and clinical perspective on emerging arboviruses. *The Lancet Infectious Diseases*, 21
703 (5), e123–e133. [https://doi.org/10.1016/S1473-3099\(20\)30491-6](https://doi.org/10.1016/S1473-3099(20)30491-6)

704 Zhang, Z., Cai, Z., Tan, Z., Lu, C., Jiang, T., Zhang, G., & Peng, Y. (2019). Rapid i
705 dentification of human-infecting viruses. *Transboundary and Emerging Diseases*, 66
706 (6), 2517–2522. <https://doi.org/10.1111/tbed.13314>

707 Zhao, L., Luo, H., Huang, D., Yu, P., Dong, Q., Mwaliko, C., Atoni, E., Nyaruaba,
708 R., Yuan, J., Zhang, G., Bente, D., Yuan, Z., & Xia, H. (2020). Pathogenesis and I
709 mmune Response of Ebinur Lake Virus: A Newly Identified Orthobunyavirus That
710 Exhibited Strong Virulence in Mice. *Frontiers in Microbiology*, 11, 625661. <https://doi.org/10.3389/fmicb.2020.625661>

712 Zhao, L., Yu, P., Shi, C., Jia, L., Evans, A., Wang, X., Wu, Q., Xiong, G., Ming, Z.,
713 Salazar, F., Agwanda, B., Bente, D., Wang, F., Liu, D., Yuan, Z., & Xia, H. (202
714 2). *Global mosquito virome profiling and mosquito spatial diffusion pathways reveal*
715 *ed by marker-viruses* [Preprint]. *Microbiology*. [https://doi.org/10.1101/2022.09.24.5093](https://doi.org/10.1101/2022.09.24.509300)
716 00

717

718

719 **Table 1: Summary of Epidemiological characteristics in regression model.** A detailed
 720 summary of the 37 epidemiological characteristics considered in our regression model.

Classification	Name of characteristics	Detailed description of characteristics	
virus	Virus_Group	Arboviruses or ISVs	
	virus	virus name	
	vi_G	virus genus	
	vi_F	virus family	
	Count	virus counts	
vector hosts	baltimore	virus baltimore classification	
	vector_G	vector genus	
	vector_F	vector family	
	vector_O	vector order	
	vector_C	vector class	
	continent	vector continent	
	country	vector country	
	climate	vector climate	
	cross_vector_G	Counts of cross-vector host genera	
	cross_vector_F	Counts of cross-vector host families	
	cross_vector_O	Counts of cross-vector host orders	
	cross_vector_C	Counts of cross-vector host classes	
	total_vector	Total counts of cross-vector hosts	
	vector_G_T	Total counts of cross-vector host genera	
	vector_F_T	Total counts of cross-vector host families	
	vector_O_T	Total counts of cross-vector host orders	
	vector_C_T	Total counts of cross-vector host classes	
	non-vector hosts	Aves	Aves host
		Carnivora	Carnivora host
		Rodentia	Rodentia host
Chiroptera		Chiroptera host	
Primates		Primates host	
homo		homo host	
Didelphimorphia		Didelphimorphia host	
Artiodactyla		Artiodactyla host	
Perissodactyla		Perissodactyla host	
Eulipotyphla		Eulipotyphla host	
Reptilia		Reptilia host	
Lagomorpha		Lagomorpha host	
Anura		Anura host	
Pilosa		Pilosa host	
Diprotodontia		Diprotodontia host	
cross_host	total cross non-vectors host		
homo_infected	Whether or not homuns are infectious		

721 **Table 2: Summary of FunSoCs annotation results from SeqScreen.** Counts and definitions
722 of 10 distinct FunSoCs identified in this dataset.

FunSoC title	Counts	FunSoC definition
Viral adhesion	44482	Mediates viral adherence to host cells
Viral counter signaling	35256	Viral suppression of host immune signaling within host cells to avoid inflammatory responses
Host xenophagy	33656	Target host xenophagy/autophagy
Viral invasion	4376	Mediates viral invasion into host cell
Host transcription	949	Target host transcription to inhibit or activate
Host ubiquitin	880	Target host ubiquitination machinery
Host cell death	802	Target host apoptotic cell death pathways either to inhibit or activate
Resist complement	144	Enable resistance from host complement components
Antibiotic resistance	4	Counters the effect of antibiotics administered to inhibit the growth or vital functioning of bacterial or eukaryotic parasites.
Induce inflammation	1	Directly activate host inflammatory pathways to cause damage

723

724

725

726 **Figure 1: Hematophagous arthropod vector and non-vector hosts characteristics in the**
727 **dataset.** (A) The global distribution and quantity of blood-sucking vectors and their carriers.
728 (B) The number of vector hosts, the continents where they are located, and the types of hostile
729 weather conditions where they are found. (C) The characteristics of the number of non-vector
730 hosts. (D) Viruses transmitted across non-vector hosts Quantity. The abscissa is the number
731 across non-vector hosts, and the ordinate is the total number of viruses.
732
733

734 **Figure 2: Relative importance of different macroscopic characteristics in the regression**

735 **model.** The weight contributions of diverse epidemiological features of viruses in the

736 regression model to human pathogenicity.

737

738

739 **Figure 3: Hosts distribution in viral functions annotation.** The distribution of hosts for
740 known non-pathogenic viruses to humans (A) and known human-pathogenic viruses (B). The
741 actual counts of viruses are converted to percentage representations in their respective
742 sections of the chart.
743
744

745 **Figure 4: Metrics of comprehensive assessment of model performance.** The utilization of
746 ROC Curve (A)and Confusion Matrix (B) for assessing the performance of the model.
747
748

749 **Figure 5: Ranking of metrics presented in the XGBoost functional annotations model.**

750 Within the results of the XGBoost model, the functional feature importance outcomes of gain

751 (A), cover (B), and weight (C) are separately obtained. These three metrics collectively reflect

752 the relative significance in determining the pathogenicity of hematophagous arthropod

753 vector-borne viruses.

754

755

756 **Figure 6: The collective impact of viral function annotations on pathogenicity prediction**
757 **analyzed through SHAP.** Providing a comprehensive overview of how various viral function
758 annotations collectively contribute to the model's predictions regarding pathogenicity.
759
760

761 **Figure 7: Detailed analysis of the interactions among crucial features in pathogenicity**
762 **prediction models through SHAP.** The interactions examined include those between viral
763 sequence size and viral adhesion (A), host xenophagy and viral adhesion (B), as well as viral
764 invasion and viral counter signaling (C). These analyses contribute to a deeper understanding
765 of the combined influence of these features on pathogenicity predictions.
766
767
768

769 **Supplementary Information**

770 **Table supplement 1: Family for hematophagous arthropod vector**

Family
Culicidae
Phlebotominae
Ceratopogonidae
Simuliidae
Tabanidae
Cimicidae
Ixodidae
Argasidae
Stenoponiidae
Phthiraptera

771

772

773 **Table supplement 2: Hyperparameter settings for the XGBoost regression model.**

774 Optimized parameter settings for the XGBoost regression model obtained through rigorous
775 experimentation and fine-tuning.

Hyperparameter	Value
booster	dart
eta	0.15
max_depth	3
subsample	0.7
objective	reg:logistic
tree_method	exact
max_cat_threshold	20
eval_metric	["logloss", "auc", 'error']

776

777

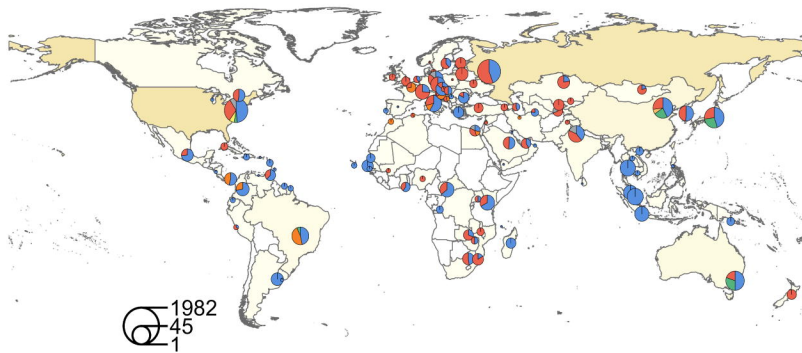
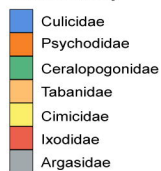
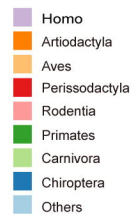
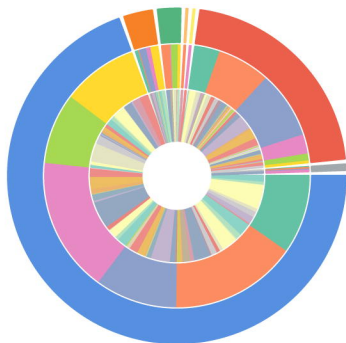
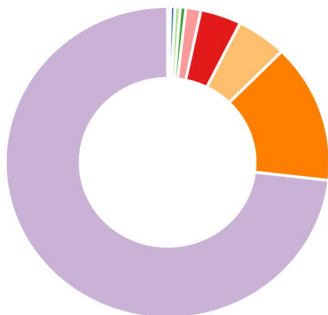
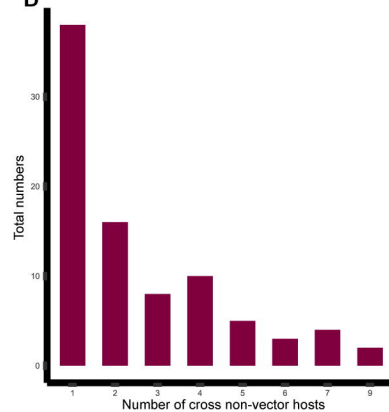
778

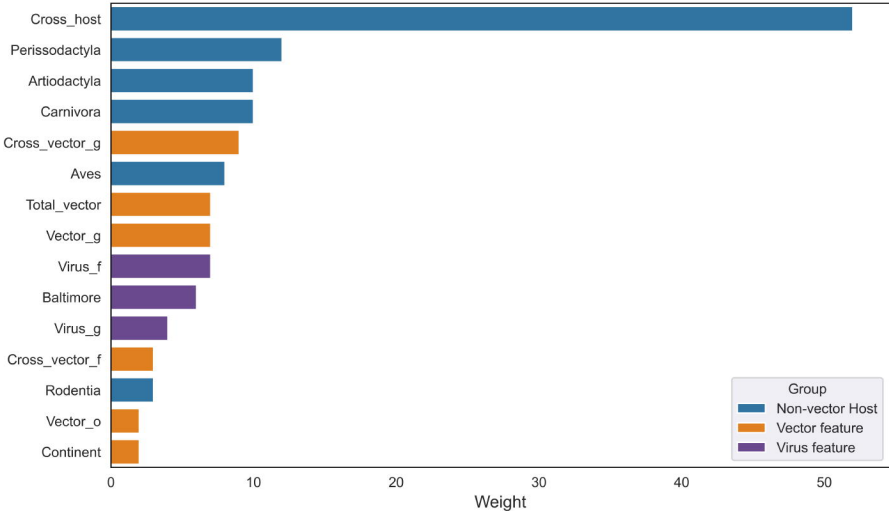
779

780 **Table supplement 3: Hyperparameter settings for the XGBoost classification model.**
781 Optimized parameter settings for the XGBoost classification model obtained through rigorous
782 experimentation and fine-tuning.

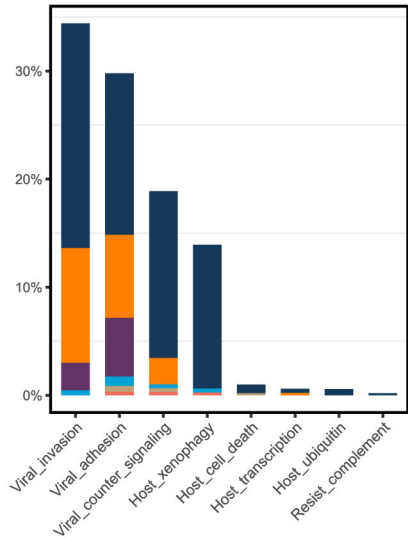
Hyperparameter	Value
objective	binary:logistic
tree_method	exact
scale_pos_weight	0.26
eta	0.15

783

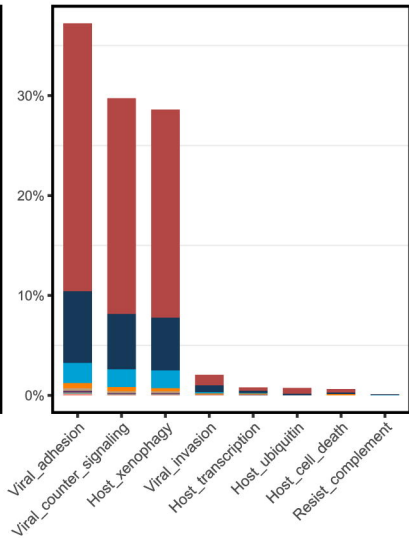
A**Numbers of viruses****Continent****Climate****Vectors Family****Non-vector hosts****B****C****D**



A

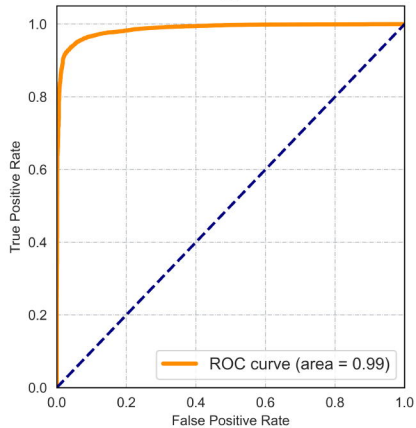
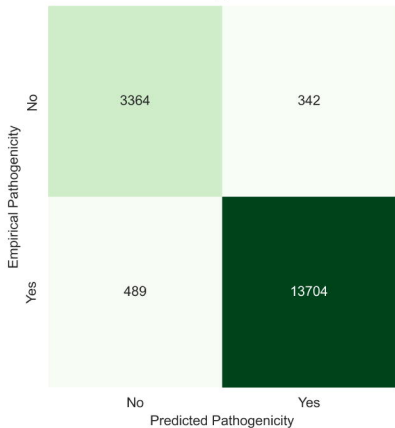


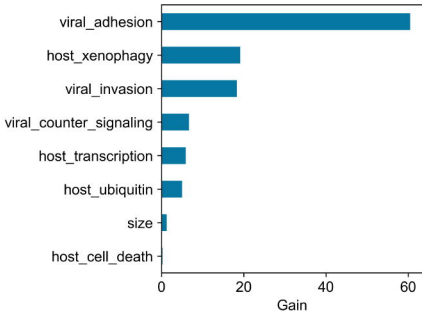
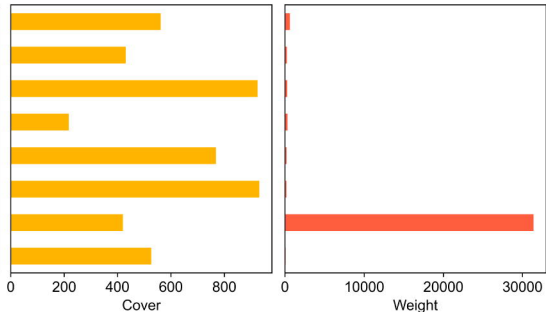
B



Category

- Homo
- Vectors
- Aves
- Artiodactyla
- Rodentia
- Perissodactyla
- Primates
- Carnivora
- Chiroptera

A**B**

A**B****C**