1       # Viral genomic features predict *Orthopoxvirus* reservoir hosts

2    Katie K. Tseng[1], Heather Koehler[2], Daniel J. Becker[3], Rory Gibb[4,5], Colin J. Carlson[6], Maria del

3    Pilar Fernandez[1*¶], and Stephanie N. Seifert[1*¶]

4

5

6

7    [1] Paul G. Allen School for Global Health, Washington State University, Pullman, Washington,
8    United States of America
9

10    [2] School of Molecular Biosciences, Washington State University, Pullman, Washington, United
11    States of America
12

13    [3] School of Biological Sciences, University of Oklahoma, Norman, Oklahoma, United States of
14    America
15

16    [4] Centre for Biodiversity and Environment Research, Department of Genetics, Evolution and
17    Environment, University College London, London, United Kingdom
18

19    [5] People & Nature Lab, UCL East, University College London, Stratford, London, United Kindom
20

21    [6] Center for Global Health Science and Security, Georgetown University, Washington, DC,
22    United States of America
23

24

25    * Corresponding authors: Stephanie N. Seifert and M. Pilar Fernandez.

26    Email: stephanie.seifert@wsu.edu (SNS) and pilar.fernandez@wsu.edu (MPF)

27

28

29    ¶ These authors jointly supervised this work.

30  **Author Contributions:** K.K.T, P.F. and S.N.S. designed research; K.K.T., H.K., M.P.F. and

31  S.N.S. performed research; D.J.B., R.G. and C.J.C. contributed to analytic tools; K.K.T, H.K.,

32  M.P.F. and S.N.S. analyzed data; and K.K.T., H.K., M.P.F. and S.N.S. wrote the paper. All

33  authors reviewed the paper.

34  **Competing Interest Statement:** The authors declare no competing interest.

35  **Classification:** Biological Sciences and Ecology.

36  **Keywords:** Infectious diseases | prediction | machine learning | boosted regression trees |

37  orthopoxviruses, Mpox

38

## Abstract

Orthopoxviruses (OPVs), including the causative agents of smallpox and mpox have led to devastating outbreaks in human populations worldwide. However, the discontinuation of smallpox vaccination, which also provides cross-protection against related OPVs, has diminished global immunity to OPVs more broadly. We apply machine learning models incorporating both host ecological and viral genomic features to predict likely reservoirs of OPVs. We demonstrate that incorporating viral genomic features in addition to host ecological traits enhanced the accuracy of potential OPV host predictions, highlighting the importance of host-virus molecular interactions in predicting potential host species. We identify hotspots for geographic regions rich with potential OPV hosts in parts of southeast Asia, equatorial Africa, and the Amazon, revealing high overlap between regions predicted to have a high number of potential OPV host species and those with the lowest smallpox vaccination coverage, indicating a heightened risk for the emergence or establishment of zoonotic OPVs. Our findings can be used to target wildlife surveillance, particularly related to concerns about mpox establishment beyond its historical range.

## Introduction

*Variola virus*, the smallpox-causing agent belonging to the *Orthopoxvirus* genus (OPV), has left an indelible mark on human history. This exceptionally virulent disease has triggered some of the most catastrophic outbreaks in human memory, leading to significant morbidity and mortality on a global scale. However, it also catalyzed progress in therapeutic intervention, inspiring the development of the first and highly effective vaccine. This vaccine leveraged cross-protective immunity from a closely related but less virulent OPV, ultimately leading to the successful eradication of smallpox[1]. In 1980, smallpox became the first human disease to be eradicated as a result of a global, coordinated vaccination effort --- and only made possible by a lack of suitable animal reservoirs to maintain variola virus outside of human populations[2]. After the eradication of smallpox, vaccinations were largely discontinued, resulting in a worldwide reduction in immunological memory against OPVs[3]. Nevertheless, the OPV genus is a diverse group, with many members still circulating in animal reservoirs, facilitating periodic emergence in humans and complicating intervention efforts.

3

69      Orthopoxviruses are notable for their varied mammalian host breadth and pathogenicity,

70      although the full range of circulation in wildlife is unknown [4]. The diversity in host breath has

71      been attributed to variations in the large repertoire of accessory genes found across poxvirus

72      species including many genes associated with the inhibition of host innate immune responses

73      [6–8]. The evolutionary dynamics of OPVs is marked by gains and losses among accessory

74      genes, with genome reduction thought to contribute to host specialization [4]. For example, the

75      modified vaccinia virus Ankara, a third-generation smallpox vaccine, lost considerable genetic

76      information (roughly 15% of the ~200kb genome), including many genes used by OPVs to

77      regulate the mammalian host environment, resulting in a restricted host range [9]. The viral

78      genome of cetaceanpox viruses, a putative sister clade to OPVs, encodes for roughly half the

79      number of proteins found in other poxviruses (7), and these viruses are thought to have highly

80      restricted host ranges. Other OPVs, including mpox virus (formerly known as monkeypox virus)

81      and cowpox virus, are capable of infecting a broad range of hosts, increasing the likelihood of

82      recurrent spillover events into the human population [4]. The unprecedented recent global

83      spread of mpox virus has raised concerns about human-to-animal spillback to susceptible hosts

84      outside its historical range, as has been observed with SARS-CoV-2 in white-tailed deer

85      populations [11]. Furthermore, the recent emergence of the novel borealpox virus (formerly

86      known as alaskapox virus) in humans from an unknown animal reservoir suggests a high

87      likelihood of unidentified OPVs in nature with zoonotic potential [12].

88      Given this complex landscape of host ranges and zoonotic potential, identifying possible

89      animal reservoir hosts of OPVs is critical for anticipating spillover events. However, identifying

90      reservoirs through field sampling is a resource-intensive endeavor [13]. To target these efforts,

91      statistical models can be used to predict undiscovered hosts of zoonotic viruses. These

92      statistical models leverage intrinsic traits of the virus [14], host [15,16], or both to predict

93      "missing links" in the host-virus network (*link prediction*), or can focus more specifically on host

94      susceptibility to a specific pathogen or pathogen group of interest *(host prediction)*. This latter

95      approach is more easily implemented, given easy access to datasets that capture host

96      ecological, geographic, morphological, and phylogenetic features [16], and a lack of comparable

97      standardized datasets of "viral traits."

98      These types of modeling studies may be undermined by their reliance on host ecological

99      and phylogenetic features, and conversely, the lack of predictors that capture host-virus

100     compatibility at molecular scales. For example, several such modeling efforts have indicated a

101     high probability of compatibility between domestic pigs (*Sus scrofa*) and SARS-CoV-2

102     [16,18,19], which was further supported by *in vitro* susceptibility of pig-derived cell lines [20].

103   However, *in vivo* challenge studies demonstrated that SARS-CoV-2 fails to infect domestic pigs

104   [21,22], suggesting a post-entry incompatibility that is poorly understood or captured in current

105   predictive models. Similarly, the Egyptian fruit bat, *Rousettus aegyptiacus*, was predicted to be

106   a compatible host for Nipah virus using host trait data [23], but failed to support Nipah virus

107   replication with experimental challenge [24].  Furthermore, viruses evolve over time, changing

108   host breadth or transmission potential, as illustrated by the expanded host range of the SARS-

109   CoV-2 Omicron variant [25]. As the evolution of viruses is shaped by the host environment, from

110   cell entry to evading host innate and adaptive immune response, viral genomes encode

111   valuable information for predicting their host range but have seldom been included in predictive

112   modeling of host-virus associations.

113         Here, we developed a trait-based approach using boosted regression trees (BRTs), a

114   machine learning algorithm commonly used in ecological and evolutionary research, and

115   integrated both host ecological traits and viral genomic features to predict mammal-OPV

116   associations. To provide a basis for comparison, we constructed traditional host prediction

117   models that used only host traits to predict OPV positivity across mammal genera. These

118   models, trained separately on datasets of candidate host genera with evidence of exposure to

119   OPV species (based on molecular detection of viral genetic material) or likely susceptible to

120   infection (based on isolation of an OPV from a host), assess how different detection methods

121   influence predictions. Building on this framework, we combined PCR and virus isolation data

122   with OPV whole genome sequences to create a comprehensive dataset of host-virus

123   interactions. Using this enriched dataset, we developed a model to predict the probability of a

124   specific link between a single OPV and a single host genus (i.e., *link prediction*) by including

125   both host traits and virus features critical to host-virus interactions on a molecular scale. Unlike

126   a host prediction model, which predicts positivity for a virus or group of viruses, the link

127   prediction model predicts compatible host-virus pairs, an approach that can be applied to

128   optimize targeted sampling and provide quantitative insights into the hypothesized role of viral

129   genomic variation in shaping OPV host range.

# Results

## Host prediction models

132   Host prediction BRT models aimed at predicting the host range of OPVs using either known

133   host exposure to OPVs (based on PCR) or known susceptible hosts of OPVs (based on virus

134    isolation) had moderate predictive accuracy (Fig 1; S1 Table). Compared to the host exposure
135    model, the susceptible host model had higher overall model accuracy (area under the receiver
136    operating characteristic curve [AUC] = 0.88 vs. 0.86; $t$ = -4.38, $p<0.001$), higher specificity (t = -
137    7.21, $p < 0.001$), and lower sensitivity ($t$ = 5.66, $p < 0.001$) (Fig 1; S1 Table). Likewise, the
138    predicted probabilities of OPV positivity were significantly correlated between the two evidence
139    type models (Spearman ρ = 0.535, $p<0.001$), with predictions of the host exposure model and
140    the susceptible host model both exhibiting strong phylogenetic signals (Pagel's λ = 0.80 and
141    0.90, respectively). Phylogenetic factorization identified some overlapping taxonomic patterns in
142    the predictions of both host prediction models, including a higher mean probability of genera
143    from the family Felidae ($n$ = 17 species) to host OPVs, and a lower mean probability of the
144    orders Lagomorpha and Rodentia ($n$ = 514) as well as a subclade of the rodent suborder
145    Hystricomorpha ($n$ = 74) to host OPVs (Figs 2 and S1; S2 and S3 Tables).
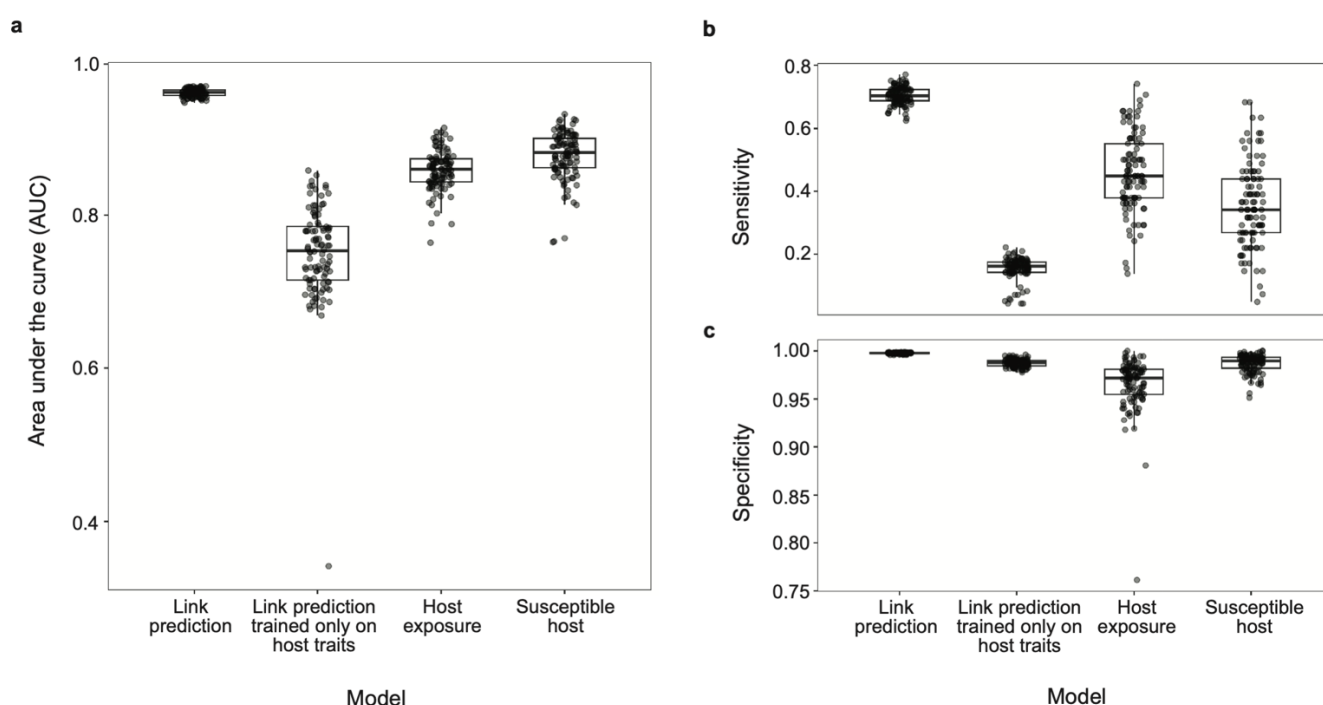
146



147
148    **Fig 1. Performance measures of boosted regression tree models.** Model accuracy was
149    evaluated by (a) area under the receiver operating characteristic curve, (b) sensitivity and (c)
150    specificity. Link prediction models were trained on known host-virus links as the response variable
151    but differed in the inclusion versus exclusion of viral genomic traits as predictors. Host prediction
152    models were trained on RT-PCR positivity (i.e., host exposure model) versus virus isolation data
153    (i.e., susceptible host model). All models were trained on 100 random splits of training (70%) and
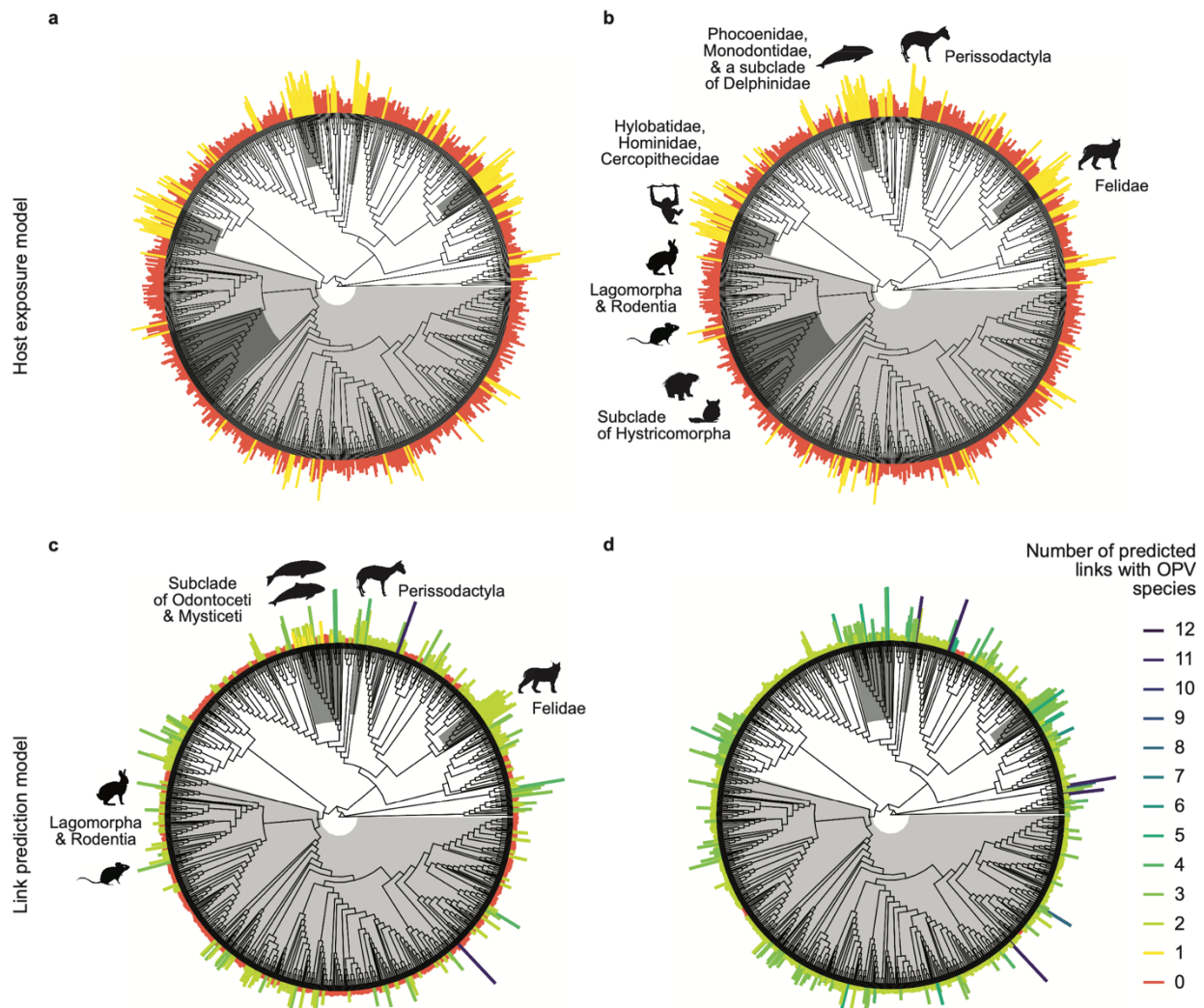154    test (30%) data.
155
156

6

**Fig 2. Predictions of *Orthopoxvirus* positivity reveal taxonomic patterns and the effects of threshold moving.** For (a, b) the host exposure model and (c, d) the link prediction model, bars/segments are scaled by predicted probabilities and colored by binary classification of predicted host genera assuming an 80% sensitivity threshold (a, c) or a threshold maximizing the sum of sensitivity and specificity (b, d). For the link prediction model, if multiple links were predicted per host genus, their predicted probabilities were averaged to represent the genus' mean probability of hosting any OPV. Clades identified through phylogenetic factorization with significantly different mean predictions are shaded in grey with corresponding labels and representative mammal silhouettes, which apply across panels of the same model: (a, b) for host exposure and (c, d) for link prediction

## Link prediction model

In contrast, the link prediction model trained on host traits and viral features (to predict the existence of a host-virus link) classified compatible host-virus pairs with higher accuracy, higher specificity, and higher, moderate sensitivity (Fig 1; S1 Table). To explore the robustness of the link prediction model, we tried predicting links solely on host traits, which, similar to host prediction models, resulted in lower model accuracy, specificity and sensitivity compared to

7

175   training on both host traits and viral features (Fig 1; S1 Table). We also tried excluding host
176   associations with vaccinia virus from our link prediction model trained on host and virus traits.
177   This decision was based on the fact that multiple passages *in vitro* of vaccinia virus have led to
178   the artificial selection of viral genes implicated in virulence, replication capacity, and host range,
179   which may not reflect natural host-virus interactions. Excluding links with vaccinia virus
180   significantly increased sensitivity ($t$ = -10.06, $p$ < 0.001) and lowered specificity ($t$ = 2.99, $p$ <
181   0.001) and overall model accuracy, albeit marginally (AUC = 0.95 vs. 0.96; $t$ = 10.1, $p$ < 0.001)
182   (Fig 1; S1 Table).
183
184   Link predicted probabilities, represented as the mean probability of hosting any OPV per host
185   genera, displayed strong phylogenetic signal (Pagel's $\lambda$ = 0.87). Similar to the host prediction
186   models, phylogenetic factorization predicted Felidae ($n$ = 17 tips) more likely to host OPV and
187   Lagomorpha and Rodentia ($n$ = 514) less likely to host OPV (Fig 2; S2-S4 Tables). The order
188   Perissodactyla ($n$ = 8) was also predicted more likely to host OPVs in agreement with the host
189   exposure model (S2-4 Tables), while overlapping taxa of the order Cetacea and Ziphiidae were
190   predicted more likely to host OPVs by both the link prediction and susceptible host models (Figs
191   2 and S1; S3 and S4 Tables).
192

# Optimizing classification

194   To address the problem of class imbalance (i.e., the unbalanced representation of hosts and
195   non-hosts or links and non-links in our training datasets), we explored various optimal
196   thresholding methods for generating binary host classifications from the predictions of our host
197   and link prediction models. Unsurprisingly, the number of predicted mammalian hosts for OPVs
198   was highly sensitive to the classification threshold, particularly for the link prediction model (Fig
199   2; S5 Table). Applying an 80% sensitivity threshold ($t_{ReqSens80}$ = 0.36) versus a 90% sensitivity
200   threshold ($t_{ReqSens90}$ = 0.28) to host exposure model predictions resulted in a 2.3-fold increase in
201   the number of predicted host genera (from $n$ = 116 to 263) (S5 Table). The same thresholds
202   selection method applied to link predictions ($t_{ReqSens80}$ = 0.13 vs. $t_{ReqSens90}$ = 0.05) resulted in a
203   3.6-fold increase in the number of predicted host genera (from $n$ = 502 using 80% sensitivity to
204   n= 1,791 when using 90% sensitivity), which had a similar effect to applying the threshold where
205   sensitivity equals specificity ($t_{Sens=Spec}$ = 0.04; $n$ = 1,864) (S5 Table). As evident across all
206   models, though, lowering the threshold value maximized sensitivity, ensuring that fewer positive
207   cases were missed, albeit at the expense of a higher false positive rate.
208   Threshold selection also led to changes in the taxonomic composition of previously unobserved,
209   predicted mammal genera when grouped by mammal order. For the link prediction model,
210   increasing the required sensitivity from 80% to 90%, which lowered the threshold value, led to
211   an increased representation of rodents (from 34% to 42%) and a decreased representation of
212   carnivores (from 23% to 17%), while the percentage of unobserved predicted genera from other
213   mammal orders remained approximately the same: e.g., Cetartiodactyla (15-16%), primates
214   (10-11%), and Eulipotyphla (7-9%).

## Feature importance

We observed some overlap in the mammal traits identified as predictive of OPV positivity and compatible host-virus pairs. Across host trait models and both link prediction models (trained on host and viral traits vs. host traits only), PubMed citations of host genera, as an indicator of sampling effort, and dispersal potential (i.e., the distance an animal can travel between its place of birth and its place of reproduction) were consistently important predictors (Fig 3; S2 and S3 Figs; S6 and S7 Tables). To determine whether host citation counts, as a proxy for sampling effort, confounds the relationship between host trait profiles and OPV positivity, we conducted a secondary set of BRTs that modeled citation counts as a Poisson response, and found that host traits did not predict study effort based on its classification performance (mean AUC = 0.5, which is not different from chance). The two host traits with the highest relative importance in the in the model including both the viral and host traits (Fig 3C), which had the best classification performance, were *island dwelling* (the proportion of species in a genus having 20% or more of the breeding range occurs on an island) and *dispersal* (median distance in km travelled by species in a genus, between the birth site and the breeding site). Both variables followed a positive association with the probability of hosting OPVs, although their effect sizes were small (S4 Fig).
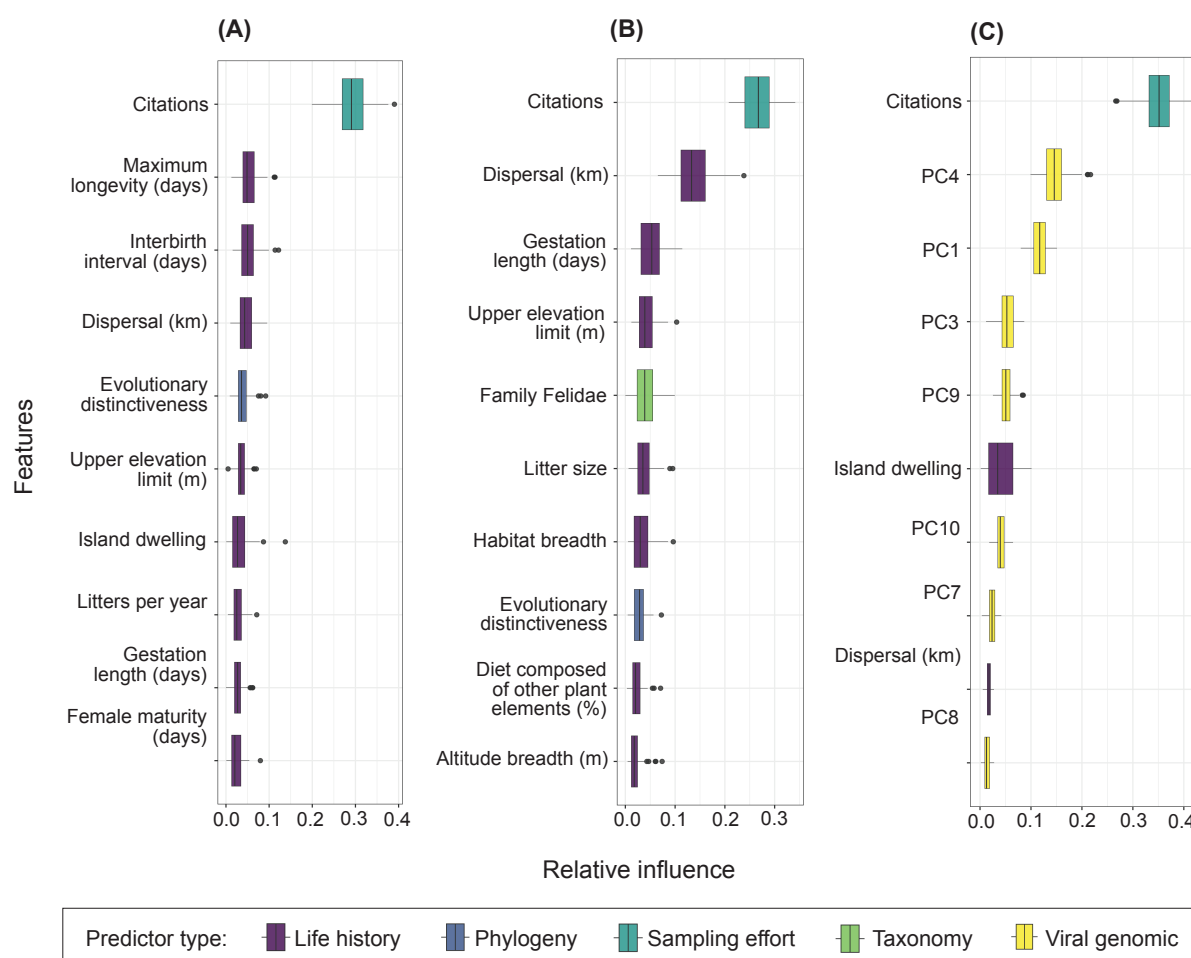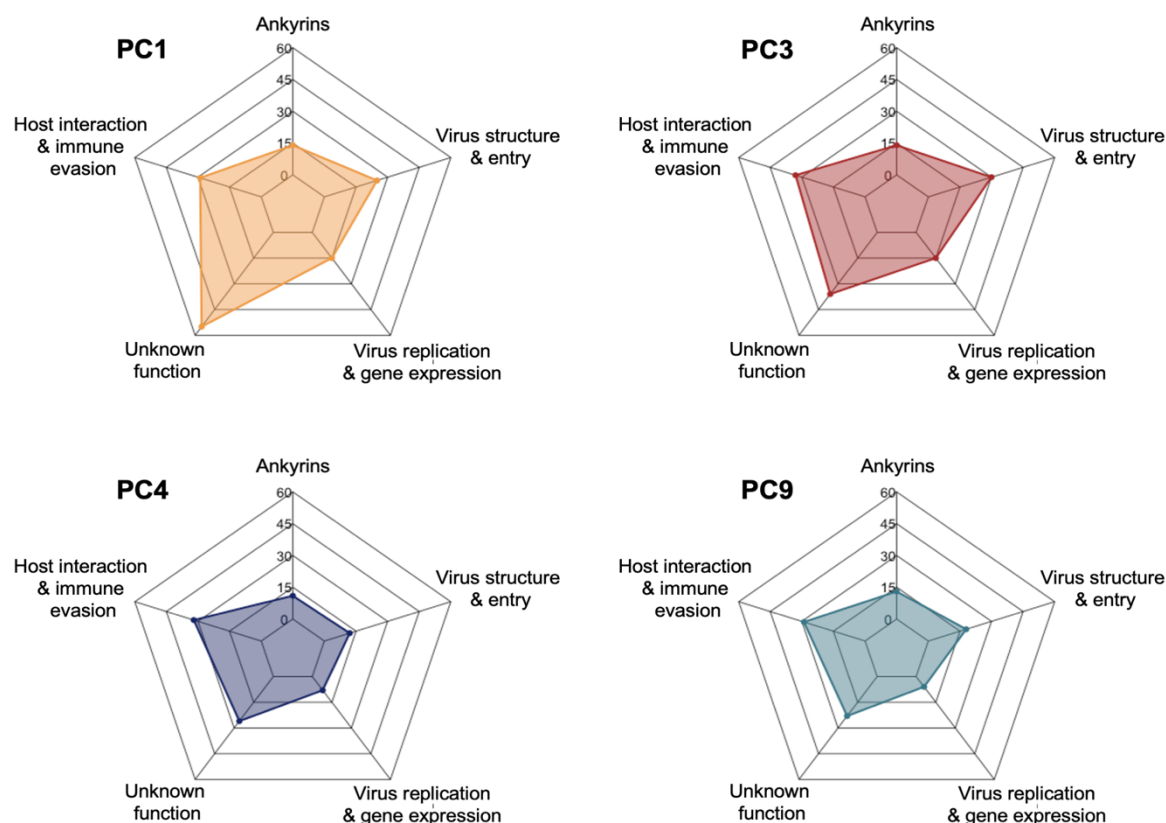
234  **Fig 3. Relative influence of model features ranked.** The most important predictor variables in
235  BRT analysis are shown for (a) the host exposure model, (b) the susceptible host model, and (c)
236  the link prediction model trained on a combination of host-virus traits, including the principal
237  components (PCs) of viral accessory gene data, with explicit pairing of host-virus links as the
238  response. Each horizontal bar plots the variability in the relative influence of the predictor
239  variables as measured across 100 random partitions of training (70%) and test (30%) data.
240  Boxplots show the median and interquartile range, whiskers are the extremes, and circles are
241  additional outliers.
242
243  In our link prediction model trained on host and viral traits, viral features were represented by
244  ten principal components (PCs) derived from a principal component analysis (PCA; Supporting
245  Methods). These components captured the majority (70%) of the variance in the
246  presence/absence of viral accessory genes (S5 Fig). While both host prediction models and the
247  link prediction model trained solely on host traits ranked life history traits highly (Figs 3A and 3B;
248  S6 Fig), rankings of relative feature importance for the link prediction model identified a
249  multitude of PC variables as important predictors (Fig 3C). Further analysis of the influential
250  PCs and their loadings, which indicate the relative contribution of individual genes to each PC,
251  found potential patterns in the functional roles of genes with high positive or negative loadings
252  (S1 Data). For instance, genes with the greatest contributions to PC4 and PC3 largely encoded
253  proteins involved in host interaction and immune evasion, including non-essential membrane
254  proteins and proteins that target innate cytokines/chemokine and cell death inhibitors (Fig 4; S7
255  Fig). Notably, for PC1, genes with the greatest positive loadings were predominately associated
256  with cowpox virus, while genes with the greatest negative loadings were primarily associated
257  with mpox virus (S1 Data), indicating clustering among specific OPV species (S8 Fig) that may
258  not have been fully captured by the accessory genes included in the PCA. A list of NCBI
259  accession numbers and PCA loadings can be found in the github repository
260  (github.com/viralemergence/PoxHost/data).
261

262



263
**Fig 4. Relative contribution of predicted functional groupings among genes influential to principal components (PC) 1, 3, 4, and 9.** For the most influential viral traits based on link prediction, we assigned predictive functions to genes with loading values greater than 1.5 times the standard deviation in the positive range or less than 1.5 times the standard deviation in the negative range ($n$ = 89,128, 123, and 89, respectively). We then plotted counts for each gene category (ankyrins, virus structure and entry, virus replication and gene expression, host interactions and immune evasion or unknown function) for PC1, PC3, PC4, and PC9.

## Distribution of potential *Orthopoxvirus* hosts

Mapping of the geographic distribution of observed hosts (i.e., known links) alongside predicted host genera (observed and unobserved) revealed novel hotspots of overlapping OPV hosts in parts of Indonesia and Malaysia, southern East Africa, the West African coastline, the Amazon basin, and the Brazilian coastline (Figs 5A and 5B). We identified the highest concentration of potential OPV hosts in the Eastern Himalayas, western Cameroon, and regions of Central and Eastern Africa extending from the Democratic Republic of Congo to Kenya (Fig 5B). Similar comparisons of observed vs. predicted host genera for mpox virus revealed novel hotspots in the rainforested regions of Central and Eastern Africa and parts of the southern and northern US plains (S9A and S9B Fig). Our models show high overlap between geographic regions with the most potential OPV host species (Fig 5B) and areas where the smallest percentages of the population have been vaccinated against smallpox [26](Fig 5C) and are, therefore, at higher risk of zoonotic OPV emergence.
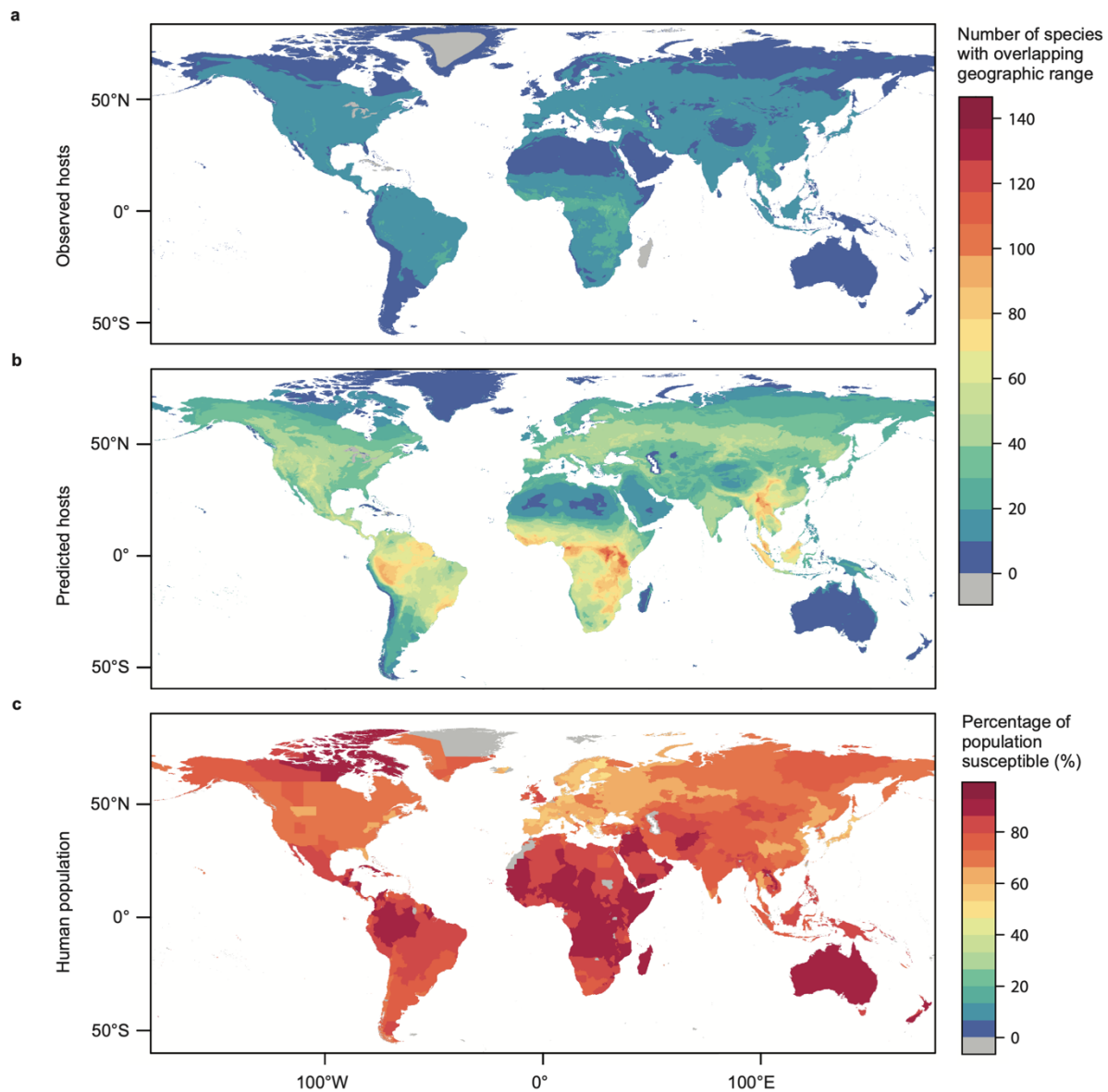
285



286

**Fig 5. Geographic distribution of *Orthopoxvirus* hosts.** Host distributions are based on the IUCN Red List database of mammal geographic ranges for those species belonging to (A) observed host genera and (a) predicted (observed and unobserved) host genera based on the results of the link prediction model and applying a 90% sensitivity threshold for host classification. The corresponding legend depicts the number of species with overlapping geographic range by color. (c) depicts smallpox vaccination coverage among humans (adapted from Taube *et al.* 2022), whereby the legend indicates the percentage of the population susceptible to orthopoxviruses.

12

# Discussion

Machine learning models that leverage ecology to predict likely reservoirs of emerging pathogens can be remarkably accurate and can be used to guide the selection of target host species for monitoring and surveillance. However, trait selection based on host ecology alone can overlook the potential influence of genetic diversity on the ability of pathogens to infect a broad or narrow range of hosts. Here, we focused our study on OPVs, a group of closely related viruses whose genetic variability is thought to contribute to their diverse host range. Using BRTs trained on observed associations between OPVs and mammal hosts, we found that a model informed by both viral genomic features and host ecological traits predicted mammal hosts with improved accuracy compared to a model trained on host traits alone. Moreover, variables capturing viral accessory gene data were found to be highly predictive of host-virus compatibility, illustrating how integrating viral genomic traits, such as the presence or absence of virulence genes, can improve host prediction and help identify key features associated with host specificity. Relating the most influential PCs in model prediction to their original features, we classified genes with the highest loadings by their predicted functional roles, which offers new opportunities to explore the roles of accessory genes in the prediction of OPV host range.

Genomic features such as phylogenetic relatedness and codon usage bias can be informative of host prediction [14,27,28] but are rarely used alongside host ecological traits to inform model prediction. Previous work has demonstrated the utility of ensemble modeling for a multi-perspective, host-virus-network framework in predicting viral host ranges by integrating multiple similarity learners derived from three perspectives: genomic (i.e., nucleotide bias, codon usage, secondary structure features, and genomic dissimilarity), ecological (mammalian traits), and network-based [29]. The network perspective utilized a bipartite graph of known virus-host interactions to quantify node similarities, enabling predictions for unobserved virus-host associations based on how viruses and hosts are positioned in the network. All three perspectives were combined via ensemble modeling to enhance predicting accuracy, but were originally modeled separately. Building on this foundation, we integrate genomic, ecological, and network-based perspectives into a unified model to address the underestimation of host-virus associations. This approach demonstrates how incorporating viral genomic features in link prediction can also serve as a valuable tool for public health surveillance and generate new hypotheses about host-virus compatibility. Our study found that BRTs trained on both host and virus traits produced similar and on average higher AUC, specificity, and sensitivity than BRTs trained solely on host traits alone. The architecture of our model was qualitatively similar to that of Blagrove et al. [17], which used a link prediction approach incorporating both host and viral traits to predict potential hosts of poxviruses. Despite differences in feature selection and methodology, their model, like ours, had low to moderate sensitivity, which was attributed to an imbalance in the distribution of observed versus unobserved host-virus associations. While Blagrove et al. corrected for class imbalance using over-sampling methods, we opted to explore the effects of probability threshold selection on host predictions, offering an alternative and more informative approach to addressing this challenge..

We note that in both our host and link prediction models, citation counts of host genera was highly ranked (Fig 3), reflecting the importance of sampling effort in accurate prediction, as genera with higher citation counts are more likely to have been studied for viral associations and have greater representation among OPV sequence data. A dataset rich in variation in both viral genomic features and host species improves generalizability and out-of-sample prediction of the model. We note that the family Felidae is well-represented as hosts across diverse cowpox virus

13

343     sequences which may contribute to the high predicted probability of detecting OPVs in the
344     family Felidae. Nearly half of all represented OPV genomes originated in humans. In this
345     dataset, mpox virus, cowpox virus, and vaccinia virus are relatively well-represented while
346     borealpox virus and cetaceanpox virus are the least represented with a single host and viral
347     genome. We included citation counts as an explanatory variable to take into consideration the
348     impact of imbalanced sampling effort across taxa.
349
350     The problem of class imbalance is a common challenge when predicting the host range of
351     emerging pathogens. Not only is there limited data on rare or novel hosts, but also natural
352     variation in species abundance and host susceptibility along with sampling bias can lead to
353     imbalance in the representation of hosts and non-hosts (i.e., compatible and non-compatible
354     host-virus pairs). When class imbalance is severe, applying a default threshold of 0.5 will often
355     result in poor model performance. Thus, handling class imbalance is an important step in
356     predictive modeling. In our study, we elected to explore higher sensitivity thresholds (e.g., 80%
357     and 90% sensitivity) to reduce the risk of false negatives (i.e., potential hosts that are incorrectly
358     classified as non-hosts). In the context of zoonotic host prediction, false negatives can lead to
359     missed opportunities for surveillance and limit our understanding of the diversity of the zoonotic
360     risk landscape for various OPV species. Thus, this trade-off ensures that fewer true hosts are
361     overlooked, making our predictions more useful for proactive and precautionary measures in
362     public health and zoonotic disease surveillance. Albeit a simple approach, tuning of the
363     threshold parameter can have large implications on the interpretation of model predictions
364     (Table S5). For instance, our PCR host exposure model trained on host traits alone estimated a
365     2.0 to 4.5-fold increase in the number of potential OPV hosts, assuming an 80% and 90%
366     sensitivity threshold. Applying the same thresholds to our link prediction model trained on host
367     and viral traits resulted in a 5.0 to 17.7-fold increase in the number of potential hosts, a likely
368     artefact of the large proportion of pseudoabsences that construct our link prediction dataset.
369     This severe class imbalance yielded lower predicted probabilities across BRTs making our
370     ensemble of classifiers highly sensitive to threshold selection as compared to BRTs trained to
371     predict OPV positivity on host traits alone.
372
373     Given the wide range of predictions that can result from threshold selection, choosing the
374     optimal thresholding method should depend on the goal of the model and the context of the
375     predictive task. For instance, borealpox virus is a recently emergent OPV only identified in
376     humans thus far. While no formal surveillance data has been published, bulletins including
377     communications from the CDC have implicated voles as a putative reservoir. For the purpose of
378     this study, we proceeded with only the human host data. Based on a restrictive 80% sensitivity
379     threshold, our link prediction model predicted two potentially previously unobserved host
380     genera, *Ailurus* (red panda) and *Micromys* (harvest mice), neither of which have a geographic
381     range overlapping with known zoonotic events, and are thus, unlikely current reservoirs in the
382     observed virus geographic distribution. However, applying a less restrictive threshold of 90%
383     sensitivity resulted in the prediction of eight potential previously unobserved hosts for borealpox
384     virus, including the only known reservoirs of ectromelia virus, rodents in the cosmopolitan genus
385     *Mus* (mice). Notably, genomic analysis of borealpox virus suggests a history of recombination
386     with ectromelia virus, including in the putative gene encoding the A-type inclusion protein [30].
387     Deletion of the A-type inclusion protein in cowpox virus enhances viral replication in *Mus*
388     *musculus* [31], suggesting that variants or deletion mutations in this gene contribute to species
389     specificity. We suspect this recombination site drove the indicator for *Mus* species as a potential
390     borealpox virus reservoir in our model. In recapitulating the ectromelia virus recombination site

391    with borealpox virus, our findings suggest that when sample data are sparse, applying less
392    restrictive thresholds may be necessary to reveal informative, meaningful predictions by
393    identifying a broader set of potential hosts for guiding targeted surveillance of a rare or novel
394    virus. Conversely, when sample data are abundant, a more restrictive threshold may be
395    appropriate to identify a manageable number of potential hosts and prioritize specific candidates
396    for detailed investigation.
397
398    Similarly, threshold selection can also impact the taxonomic representation or composition of
399    model predictions. For instance, among known mammal genera susceptible to mpox virus, the
400    majority are from the Rodentia (37%) and Primate (37%) orders, which is consistent with
401    Blagrove et al. whose findings predicted a similar composition of mpox virus hosts (80% from
402    the Rodentia and Primate orders) [17]. However, we show in our link prediction model, that the
403    host composition of model predictions broadens when applying a less restrictive threshold.
404    When assuming an 80% sensitivity threshold, 45.2% of previously unobserved mpox virus hosts
405    were predicted to be from Rodentia, while increasing model sensitivity to 90% led to a decrease
406    in predicted rodent genera (31.4%) and a substantial increase in genera from the Carnivora
407    order (from 3.2% to 19.6 %), particularly those in the Felidae (cats), Canidae (canids),
408    Mephitidae *(*skunks), Mustelidae (mustelids), and Procyonidae (raccoons) families. Thus,
409    choosing the optimal threshold for classification may have important public health implications,
410    as it may expand the focus on species that may be overlooked when applying a more restrictive
411    threshold. In the mpox example, this is of particular interest since its recent global spread
412    outside endemic countries raises concerns of the potential for transmission from humans to
413    wildlife or spillback in nonendemic areas [32] and the establishment of new animal reservoirs.
414    We note that the genus *Rattus* was not predicted as a likely host for mpox virus in our link
415    prediction model at any of our tested thresholds which is consistent with laboratory experimental
416    data showing that rats are not susceptible to mpox virus and in contrast to predictions using
417    different approaches [17]. A table of predicted binary host-virus links can be found in the github
418    repository (github.com/viralemergence/PoxHost/figures/other/linkpred).
419
420    Our study reveals a striking correspondence between regions with high potential OPV host
421    species diversity and those where smallpox vaccination rates are lowest (Fig 5). This
422    juxtaposition suggests a pronounced vulnerability to zoonotic OPVs within these populations.
423    Hotspots identified as having high concentrations of postulated OPV hosts included several
424    distinct areas within the eastern Himalayas of South Asia bordering southwest China, Central
425    and Eastern Africa, Indonesia, and Malaysia – areas where wildlife sampling and monitoring of
426    spillover events should focus their efforts. Comparing the geographic distribution of predicted
427    hosts to known hosts also revealed potentially under sampled areas along the West African and
428    Brazilian coastline; these areas also coincide with many recent human cases of reemerging
429    zoonotic OPVs including mpox virus in Cameroon [33] and vaccinia virus in Brazil [34].
430    Interestingly, greater geographic dispersal (distance traveled between birth site and breeding
431    site) was not only an important predictor (Fig 3) but was also positively associated with an
432    increased probability of hosting OPVs (S2, S3 and S4 Figs), supporting theories of host
433    dispersal as a key contributor to the evolution of host-pathogen dynamics [35]. Specifically, as
434    individuals move to new areas to breed, they must adapt to new environmental conditions and,
435    in the process, are exposed to and may host more diverse pathogens. On the other hand, the
436    positive association between island dwelling and the probability of hosting OPVs could be
437    indicating higher prevalence or oversampling in constrained island communities, although the
438    effect size was small (S4 Fig).
439
440    Like all related studies, our study had several limitations, particularly with respect to sampling
441    bias, incomplete datasets, and model interpretability. First, in zoonotic host prediction, host-virus

442 association data is inherently incomplete, a limitation further compounded by sampling bias, as
443 host-virus interactions are disproportionately observed in well-studied taxa and regions,
444 especially for human-associated OPVs. Although we incorporated pseudoabsences to enhance
445 generalizability, our predictions may overemphasize well-documented interactions, while
446 underestimating associations involving rare or poorly sampled hosts, or novel human infecting
447 viruses. Ongoing efforts to consolidate host-virus interaction data will help to reduce biases and
448 improve data quality over time.
449
450 Second, as only complete OPV genomes were extracted from NCBI for inclusion in the link
451 prediction model, the number of host-virus associations with accessory gene data available was
452 limited. Moreover, missing data for host and virus traits precluded more precise predictions at
453 the species level. Thus, we aggregated the dataset at the genera level to reduce the number of
454 pseudoabsences in the models. We expect future prediction models to improve in their accuracy
455 as whole genome sequences of host and virus become more readily available.
456 Third, our study elected to classify accessory gene data based on presence-absence, trading
457 simplicity over potential loss of information that could have related protein activity to host-virus
458 compatibility. Subsequent studies could investigate transforming accessory gene data as
459 continuous variables based on amino acid conservation or ranked categorical variables, where
460 accessory genes are missing, truncated and non-functional, truncated but likely functional, or
461 intact, as previous studies have explored [36]. Reducing accessory gene variables to their
462 principal components (to avoid overfitting our model) also precluded us from interpreting the
463 effects of individual accessory genes in predicting host-OPV pairs. However, our findings could
464 inform future iterations by training a model on the genes identified in this study as contributing
465 highly to influential principal components, thereby directly exploring the relationship between
466 specific viral accessory genes and host range. Furthermore, many functional roles of genes
467 were predicted based on sequence homology to known OPV accessory genes, particularly for
468 vaccinia virus or mpox virus, but are not well characterized in other species. Because there is
469 no standard for naming genes, many functions are putative and gene names can be based off
470 other viruses, which can be misleading. Finally, predictive models can also be limited by the
471 machine learning algorithms they employ. As black box models, BRTs can be difficult to
472 interpret and do not provide a complete or exact understanding of the relationships and potential
473 interactions between predictors and the response variable. While our study employed cross-
474 validation to assess the model's robustness, we did not perform external validation or biological
475 validation, which requires field and laboratory studies to confirm predicted host-virus
476 interactions. Nevertheless, our model can also provide guidance for empirical sampling and
477 initial insight into the molecular signatures of host-virus compatibility. While not the aim of our
478 study, future link prediction models informed by viral genomic data could investigate training on
479 different types of infection evidence to gain molecular insights into host capacity.
480
481 Historically, zoonotic OPVs like mpox virus were associated with small, sporadic outbreaks
482 limited to endemic countries. However, large outbreaks in recent years including the worldwide
483 spread of mpox virus in 2022 warrant growing concerns about the potential for mpox virus and
484 other OPVs to establish new endemic areas. Recent evidence of animal-to-human and human-
485 to-animal transmission of mpox virus also suggest the potential for new reservoir species to
486 exist in traditionally non-endemic regions, underscoring the importance of predicting the host
487 range of emerging pathogens. Adapting models to predict compatibility is also crucial to
488 expanding the growing toolkit of analytical models used to predict host-pathogen interactions.
489 Here, we demonstrate how trait-based machine learning models can be trained on both

490  mammal traits of potential hosts and the genomic features of OPVs to improve the accuracy of
491  host prediction and gain a more comprehensive understanding of the factors that influence host-
492  virus compatibility. By exploring the trade-off between sensitivity and specificity through
493  classification threshold selection, we show the importance of aligning models to meet the
494  specific needs and objectives of host prediction for greater practical application.
495

# Materials and methods

497  We used OPVs as a case study to develop and validate a multi-perspective modeling
498  framework for the prediction of host-virus associations. We generated two types of trait-based
499  models to compare model performance and predictions: (1) a host prediction model trained
500  separately on two evidence levels for host capacity (i.e., host exposure via PCR positivity data
501  and host susceptibility via virus isolation data) using host ecological traits alone; and (2) a link
502  prediction model trained on explicit pairing of host-virus associations using a combination of
503  host and virus features. The former model type predicted host positivity for OPVs, a group of
504  viruses, while the latter predicted the existence of a host-virus link. Our modeling targets were to
505  identify taxonomic and phylogenetic patterns in the predicted reservoir hosts of OPVs along with
506  their corresponding geographic distribution, rank influential host and viral features, and
507  investigate the functional roles of accessory genes that contributed substantially to link
508  prediction.
509

## Datasets

511  For the host prediction model, we obtained known host positivity data from the Global Virome in
512  One Network (VIRION), an open database of host-virus interactions drawn from scientific
513  literature and online databases [37]. We included only host positivity detected via PCR
514  (molecular detection) or virus isolation and excluded detections with variola virus, as the virus
515  was eradicated prior to the development of many modern diagnostics. We then collapsed host
516  positivity data to the genus level and merged the data with the broader mammal taxonomy to
517  obtain pseudoabsences for mammal genera with unobserved OPV associations [38]. Next, we
518  derived host predictors from COMBINE, a published mammal dataset of 54 morphological,
519  biogeographic, and life history traits that is taxonomically compatible with the IUCN Red List v.
520  2020-2 and PHYLACINE v. 1.2 [39–41]. We applied summary measures to represent ecological
521  trait data at the genus level by transforming all categorical variables into multiple binary
522  variables. We then aggregated binary variables to the genus level assuming the mean to obtain
523  the proportion of species in a genus having the variable outcome/trait, and we aggregated
524  continuous and integer trait variables assuming the median. To account for large gaps in trait
525  coverage, we excluded variables with zero variance or with data for less than 60% of host
526  genera (S10 Fig), resulting in 61 ecological trait variables. We also generated binary predictors
527  for each mammal family to represent taxonomy, leading to an additional 42 taxonomic variables
528  in the host trait model. Lastly, we incorporated a variable for evolutionary distinctiveness
529  (*ed_equal*) using the *picante* package in R and a count variable of the number of scientific
530  publications on each host genera (*cites*) as a measure of sampling effort using the R package
531  *easyPubMed.* A complete list of host predictors incorporated in the host prediction model can be
532  found in the supporting information (S6 Table).
533

534 For the link prediction model, we merged data from VIRION with host-OPV associations
535 obtained from NCBI by extracting OPV genomes and their annotations
536 (https://www.ncbi.nlm.nih.gov/)(S2 Data). We included viral isolates with full genome sequences
537 from experimental laboratory studies when possible. Due to the abundance of available mpox
538 virus and cowpox virus genomes, we included a subset of each for all genomes with novel host
539 associations or representing divergent clades. Multiple sequences for a particular virus-host
540 combination were only included if an annotated genome was available and there were
541 differences in the presence/absence of accessory genes. Again, we merged host-virus
542 association data with the broader mammal taxonomy to accommodate out-of-sample
543 predictions, expanding our mammal-virus network to include all combinations of mammal
544 genera and OPV species. As predictors in our model, we incorporated the same host predictors
545 as in the host prediction model, with the exception of 110 taxonomic variables instead of 41.
546 Additionally, we incorporated viral predictors by identifying accessory genes from the
547 aforementioned OPV genomes using Roary (https://sanger-pathogens.github.io/Roary/) with
548 80% minimum sequence identify for blastp and 95% of isolates a gene must be in to be
549 designated as part of the core genome. After transforming the accessory gene data into a binary
550 matrix across 981 accessory genes from 197 OPV sequences, we conducted PCA to reduce
551 the dimensionality of our dataset and distill presence/absence variables down to their most
552 important features. The first ten principal components, which explained roughly 70% of data
553 variance, were then included as viral predictors in our link prediction model (S5 Fig). OPV
554 genomic data used in the PCA was available for 12 OPV species, but only for a limited number
555 of OPV-host pairings (e.g., Mpox virus-human, Mpox virus-dog). However, as we considered the
556 10 PC variables derived from the PCA as viral traits, we conducted median imputation of the PC
557 values for the links between hosts with any of those 12 OPV species for which PC data were
558 missing (e.g., for all unobserved Mpox virus-host pairings we applied the median PC value for
559 Mpox virus, for each one of the PCs). A complete list of host and viral features incorporated in
560 the link prediction model can be found in the supporting information (S7 Table).
561
562 Lastly, for phylogenetic analysis, we used a supertree of extant mammal species trimmed to the
563 genus level [38] (https://vertlife.org/data/mammals/). Data collation (cleaning and merging) were
564 conducted in R version 4.2.2. Additional details regarding the cleaning of final datasets used in
565 the host trait and link prediction models are available in the supporting information including
566 dimension reduction for generating viral predictors and taxonomic reconciliation.
567

568 ## Statistical analysis

569 ### Host prediction models
570
571 Following Mull et al. [42] (https://github.com/viralemergence/hantaro), we used a trait-based
572 model to infer the host range of OPVs, integrating only host mammal traits as predictors. This
573 approach handles binomial virus positivity of host genera as the response variable and host trait
574 data as predictors. It uses a machine learning algorithm to identify host features associated with
575 OPV detection and to make predictions of the probability of OPV detection based on how similar
576 a species' trait profile is to that of observed host species. False positive results, where a species
577 classified as a pseudoabsence is predicted to host a virus, are used to infer potential, previously
578 undetected or unobserved host genera.
579
580 To classify mammal genera as OPV hosts based on our matrix of predictors/traits, we used
581 boosted regression trees (BRTs). BRTs combine two validated machine learning algorithms:

582 regression trees, which model the relationship between outcome and predictors by recursive
583 binary splits, and boosting, an adaptive technique in which the model considers the prediction
584 errors of previous trees when fitting subsequent trees thereby improving predictive performance
585 and classification accuracy [43]. BRTs are advantageous to our study as they can handle large
586 ecological datasets, a binomial outcome variable, and a diversity of model predictors and the
587 potential interactions between them. They can also fit potentially complex nonlinear
588 relationships, produce response functions for each predictor, and are robust to missing values
589 and outliers [44].
590
591 Using BRTs, we developed separate evidence type models for classification of virus positivity: a
592 host exposure model based on PCR-detection and a susceptible host model based on virus
593 isolation data. Prior to model fitting, we used the R package *rsample* to randomly partition our
594 data into training (70%) and test (30%) sets and apply stratified sampling to ensure an equal
595 distribution of positive labels in both datasets. Next, we tuned our model by training it on
596 multiple combinations of hyperparameter values allowing for a maximum number of 5,000 trees,
597 an interaction depth ranging from 2 to 4 (to control tree complexity), a shrinkage (learning) rate
598 ranging from 0.0005 to 0.01, and a minimum number of 4 observations in terminal nodes (S11
599 Fig). Using the *gbm* package in R (Greenwell et al., 2020) for all model fitting, we applied five-
600 fold cross validation and a bag/subsampling fraction of 0.5, and we assumed a Bernoulli
601 distribution error for our binary response variable. Based on multiple model performance
602 measures, we selected the best set of parameter values for model training, which included a
603 learning rate of 0.01, and an interaction depth of three and the maximum number of trees set to
604 4500. With our tuned model for prediction, we set citation counts per genera to the mean across
605 all mammal genera to adjust for sampling effort. We then fit replicate models to 100 randomly
606 stratified partitions creating an ensemble of BRTs for each evidence type (PCR and virus
607 isolation). Additionally, we constructed a third set of BRTs to determine the influence of
608 sampling effort on the trait profiles of OPV positive genera; instead of virus positivity, we
609 modeled the citation counts of host genera as a Poisson response variable applying the same
610 optimal hyperparameter values except for 10,000 maximum trees.
611

## Link prediction model

613
614 To predict the host ranges of multiple OPV species and to allow for integration of viral genomic
615 features, we adapted the host prediction models I for link prediction whereby the model predicts
616 the existence of a host-virus link (as opposed to positivity for a virus or group of viruses). While
617 the previous host prediction models predicted OPV positivity (regardless of OPV species), the
618 link prediction model uses explicit pairing or 'links' between host genera and virus species as a
619 response variable and, therefore, allows for the inclusion of both viral and host traits
620 simultaneously as predictors. As in the host prediction models, undiscovered host-virus
621 associations are identified based on false positive predictions.
622
623 Using BRTs to predict the probability of links (associations between host genera and virus
624 species), we followed a similar workflow as the host prediction models for generating link
625 predictions. We checked that both test and training partitions contained similar proportions
626 amongst the included viral taxa. We developed three ensembles of BRTs. The first, which
627 included all combinations of host genera and OPV species (as described in *Host-virus*
628 *association data* of S1 Text), used a combination of host and viral features to predict host-virus
629 links. The second used only host features for prediction. The third used both host and viral
630 features for prediction but excluded any host-virus combination with vaccinia virus. The latter

631 was intended to explore if spontaneous, extended deletions and other modifications in viral
632 genes implicated in virulence, replication capacity and home range of vaccinia virus due to
633 multiple passages in cell cultures, may introduce bias in the genetic features included in the
634 model and affect model performance.
635

## Model performance and prediction

637
638 For each model fitting of the partitioned training and test data, we obtained model performance
639 metrics including AUC, sensitivity and specificity using the *ROCR* and *InformationValue*
640 packages in R and a threshold value of 0.5. For each ensemble model, we then calculated the
641 mean and standard error of each performance metric measured across the 100 random
642 partitions and conducted an unpaired *t*-test to compare model performance, with *p*-values
643 adjusted for the false discovery rate using the Benjamini-Hochberg correction [45,46].
644
645 To compare model predictions, we first calculated the Spearman correlation coefficient between
646 the predicted probabilities of the host exposure vs. susceptible host models as well as the link
647 prediction model when including vs. excluding vaccinia virus. Next, we assessed model
648 predictions for taxonomic patterns. Because the link prediction model predicted multiple links
649 per host genera, we calculated mean predictions per host genera to obtain a single value per
650 tree tip (branch endpoint) for phylogenetic analysis. Using the *nlme* package in R, we estimated
651 host phylogenetic signal based on Pagel's λ, a measure of the degree to which the propensity
652 for virus positivity among related genera is driven by shared evolutionary history. We then
653 identified clades with significantly different propensity for hosting OPV at various taxonomic
654 resolutions using phylogenetic factorization, a graph-partitioning algorithm that explores the
655 differences in measured traits between pairs of taxa while accounting for phylogenetic
656 dependencies [47]. To determine the significant number of phylogenetic factors (clades), we
657 adjusted for the family wise error rate using Holm's sequentially rejective approach with a 5%
658 threshold [48]. Phylogenetic factorization was implemented using the *phylofactor* package in R.
659
660 To generate binary predictions of OPV hosts and non-hosts, we transformed predicted
661 probabilities into binary classifications using the *presenceabsence* R package. Based on the
662 previously described ROC, we selected multiple thresholds to generate predictions based on
663 different methods of threshold optimization. First, we explored higher sensitivity thresholds to
664 prioritize minimizing false negatives (i.e., potential hosts that are incorrectly classified as non-
665 hosts) at the expense of increasing the number of false positives (i.e., mammal genera
666 incorrectly predicted to be hosts). We specifically selected an 80% sensitivity threshold (such
667 that 80% of observed links are detected) and a 90% sensitivity threshold, as examples of more
668 and less restrictive thresholds, or higher and lower threshold values, respectively. Second, we
669 quantified the sensitivity of results to the choice of threshold by calculating the number of
670 predicted hosts at additional threshold levels based on sensitivity-specificity trade-offs. These
671 additional approaches included finding the threshold where sensitivity equaled specificity and
672 findingthe threshold that maximized the sum of sensitivity and specificity, otherwise known as
673 the Youden Index. We then compared classifications across models and demonstrated the
674 effects of threshold selection by plotting predicted probabilities and their binary classifications on
675 a circular phylogenetic tree using the package *treeio* and *ggtree*. Lastly, we mapped the
676 geographic distribution of predicted observed and unobserved mammal genera; using the IUCN
677 Red List database of known ranges of mammal species, we layered shape files of species

678  belonging to thresholded genera to identify and contrast geographic patterns in model
679  predictions.
680
681  Finally, we ranked model features by their mean variable importance (i.e., relative influence
682  coefficients) and assessed similarities between models using the Spearman rank correlation
683  coefficient, a rank-based measure of association. Viral PC predictors with high mean variable
684  importance in our link prediction model were further analyzed by gene loading contribution to
685  identify potential patterns in the functional roles of genes with high loadings. Specifically, we
686  assigned predictive functions based on sequence homology to genes with loading values that
687  were greater than 1.5 times the standard deviation in the positive range or less than 1.5 times
688  the standard deviation in the negative range (S1 Data).
689
690  All statistical analyses were conducted in R version 4.2.2 and Microsoft Excel version 16.69.1.
691

# Data availability

693  Data used in this study is available as described in the methods section including the VIRION
694  database (https://github.com/viralemergence/virion), COMBINE: the coalesced mammal
695  database of intrinsic and extrinsic traits (Soria et al. 2021), global smallpox vaccination
696  coverage (Taube et al. 2023). Accession numbers for OPV genomes included in this study can
697  be found in Supplementary Table 2 and are accessible through GenBank
698  (www.ncbi.nlm.nih.gov).

699

# Code availability

701  The R code to reproduce the analyses is available from the GitHub repository:
702  https://github.com/viralemergence/PoxHost.

# Acknowledgments

# References

712  1.  Jacobs BL, Langland JO, Kibler KV, Denzler KL, White SD, Holechek SA, et al. Vaccinia
713      virus vaccines: Past, present and future. Antiviral Research. 2009;84: 1–13.
714      doi:10.1016/j.antiviral.2009.06.006

715   2.   Fenner F. Global Eradication of Smallpox. Reviews of Infectious Diseases. 1982;4: 916–
716        930. doi:10.1093/clinids/4.5.916

717   3.   Rimoin AW, Mulembakani PM, Johnston SC, Lloyd Smith JO, Kisalu NK, Kinkela TL, et al.
718        Major increase in human monkeypox incidence 30 years after smallpox vaccination
719        campaigns cease in the Democratic Republic of Congo. Proceedings of the National
720        Academy of Sciences. 2010;107: 16262–16267. doi:10.1073/pnas.1005769107

721   4.   Reynolds MG, Guagliardo SAJ, Nakazawa YJ, Doty JB, Mauldin MR. Understanding
722        orthopoxvirus host range and evolution: from the enigmatic to the usual suspects. Current
723        Opinion in Virology. 2018;28: 108–115. doi:10.1016/j.coviro.2017.11.012

724   5.   Sutter G, Moss B. Nonreplicating vaccinia vector efficiently expresses recombinant genes.
725        Proceedings of the National Academy of Sciences. 1992;89: 10847–10851.
726        doi:10.1073/pnas.89.22.10847

727   6.   Xiang Y, White A. Monkeypox virus emerges from the shadow of its more infamous cousin:
728        family biology matters. Emerging Microbes & Infections. 2022;11: 1768–1777.
729        doi:10.1080/22221751.2022.2095309

730   7.   Shchelkunov SN. Orthopoxvirus genes that mediate disease virulence and host tropism.
731        Adv Virol. 2012;2012: 524743. doi:10.1155/2012/524743

732   8.   Shchelkunov SN. An Increasing Danger of Zoonotic Orthopoxvirus Infections. PLOS
733        Pathogens. 2013;9: e1003756. doi:10.1371/journal.ppat.1003756

734   9.   Mayr A, Hochstein-Mintzel V, Stickl H. Abstammung, Eigenschaften und Verwendung des
735        attenuierten Vaccinia-Stammes MVA. Infection. 1975;3: 6–14. doi:10.1007/BF01641272

736   10.  Rodrigues TCS, Subramaniam K, Varsani A, McFadden G, Schaefer AM, Bossart GD, et
737        al. Genome characterization of cetaceanpox virus from a managed Indo-Pacific bottlenose
738        dolphin (Tursiops aduncus). Virus Research. 2020;278: 197861.
739        doi:10.1016/j.virusres.2020.197861

740   11.  Hale VL, Dennis PM, McBride DS, Nolting JM, Madden C, Huey D, et al. SARS-CoV-2
741        infection in free-ranging white-tailed deer. Nature. 2022;602: 481–486.
742        doi:10.1038/s41586-021-04353-x

743   12.  Springer YP, Hsu CH, Werle ZR, Olson LE, Cooper MP, Castrodale LJ, et al. Novel
744        Orthopoxvirus Infection in an Alaska Resident. Clinical Infectious Diseases. 2017;64:
745        1737–1741. doi:10.1093/cid/cix219

746   13.  Bernstein AS, Ando AW, Loch-Temzelides T, Vale MM, Li BV, Li H, et al. The costs and
747        benefits of primary prevention of zoonotic pandemics. Science Advances. 2022;8:
748        eabl4183. doi:10.1126/sciadv.abl4183

749   14.  Mollentze N, Babayan SA, Streicker DG. Identifying and prioritizing potential human-
750        infecting viruses from their genome sequences. PLOS Biology. 2021;19: e3001390.
751        doi:10.1371/journal.pbio.3001390

752    15.   Han BA, Schmidt JP, Bowden SE, Drake JM. Rodent reservoirs of future zoonotic
753         diseases. Proceedings of the National Academy of Sciences. 2015;112: 7039–7044.
754         doi:10.1073/pnas.1501598112

755    16.   Becker DJ, Albery GF, Sjodin AR, Poisot T, Bergner LM, Chen B, et al. Optimising
756         predictive models to prioritise viral discovery in zoonotic reservoirs. The Lancet Microbe.
757         2022;3: e625–e637. doi:10.1016/S2666-5247(21)00245-7

758    17.   Blagrove MS, Pilgrim J, Kotsiri A, Hui M, Baylis M, Wardeh M. Monkeypox virus shows
759         potential to infect a diverse range of native animal species across Europe, indicating high
760         risk of becoming endemic in the region. bioRxiv; 2022. p. 2022.08.13.503846.
761         doi:10.1101/2022.08.13.503846

762    18.   Fischhoff IR, Castellanos AA, Rodrigues JPGLM, Varsani A, Han BA. Predicting the
763         zoonotic capacity of mammals to transmit SARS-CoV-2. Proc Biol Sci. 2021;288:
764         20211651. doi:10.1098/rspb.2021.1651

765    19.   Wardeh M, Baylis M, Blagrove MSC. Predicting mammalian hosts in which novel
766         coronaviruses can be generated. Nat Commun. 2021;12: 780. doi:10.1038/s41467-021-
767         21034-5

768    20.   Meekins DA, Morozov I, Trujillo JD, Gaudreault NN, Bold D, Carossino M, et al.
769         Susceptibility of swine cells and domestic pigs to SARS-CoV-2. Emerging Microbes &
770         Infections. 2020;9: 2278–2288. doi:10.1080/22221751.2020.1831405

771    21.   Schlottau K, Rissmann M, Graaf A, Schön J, Sehl J, Wylezich C, et al. SARS-CoV-2 in fruit
772         bats, ferrets, pigs, and chickens: an experimental transmission study. The Lancet Microbe.
773         2020;1: e218–e225. doi:10.1016/S2666-5247(20)30089-6

774    22.   Haddock E, Callison J, Seifert SN, Okumura A, Tang-Huau T-L, Leventhal SS, et al.
775         Three-Week Old Pigs Are Not Susceptible to Productive Infection with SARS-COV-2.
776         Microorganisms. 2022;10: 407. doi:10.3390/microorganisms10020407

777    23.   Plowright RK, Becker DJ, Crowley DE, Washburne AD, Huang T, Nameer PO, et al.
778         Prioritizing surveillance of Nipah virus in India. PLOS Neglected Tropical Diseases.
779         2019;13: e0007393. doi:10.1371/journal.pntd.0007393

780    24.   Seifert SN, Letko MC, Bushmaker T, Laing ED, Saturday G, Meade-White K, et al.
781         Rousettus aegyptiacus Bats Do Not Support Productive Nipah Virus Replication. The
782         Journal of Infectious Diseases. 2020;221: S407–S413. doi:10.1093/infdis/jiz429

783    25.   Sánchez-Morales L, Sánchez-Vizcaíno JM, Pérez-Sancho M, Domínguez L, Barroso-
784         Arévalo S. The Omicron (B.1.1.529) SARS-CoV-2 variant of concern also affects
785         companion animals. Frontiers in Veterinary Science. 2022;9. Available:
786         https://www.frontiersin.org/articles/10.3389/fvets.2022.940710

787    26.   Taube JC, Rest EC, Lloyd-Smith JO, Bansal S. The global landscape of smallpox
788         vaccination history and implications for current and future orthopoxvirus susceptibility: a
789         modelling study. The Lancet Infectious Diseases. 2023;23: 454–462. doi:10.1016/S1473-
790         3099(22)00664-8

27.  Babayan SA, Orton RJ, Streicker DG. Predicting reservoir hosts and arthropod vectors from evolutionary signatures in RNA virus genomes. Science. 2018;362: 577–580. doi:10.1126/science.aap9072

28.  Brierley L, Fowler A. Predicting the animal hosts of coronaviruses from compositional biases of spike protein and whole genome sequences through machine learning. PLOS Pathogens. 2021;17: e1009149. doi:10.1371/journal.ppat.1009149

29.  Wardeh M, Blagrove MSC, Sharkey KJ, Baylis M. Divide-and-conquer: machine-learning integrates mammalian and viral traits with network features to predict virus-mammal associations. Nat Commun. 2021;12: 3954. doi:10.1038/s41467-021-24085-w

30.  Gigante CM, Gao J, Tang S, McCollum AM, Wilkins K, Reynolds MG, et al. Genome of Alaskapox Virus, a Novel Orthopoxvirus Isolated from Alaska. Viruses. 2019;11. doi:10.3390/v11080708

31.  Kastenmayer RJ, Maruri-Avidal L, Americo JL, Earl PL, Weisberg AS, Moss B. Elimination of A-type inclusion formation enhances cowpox virus replication in mice: Implications for orthopoxvirus evolution. Virology. 2014;452–453: 59–66. doi:10.1016/j.virol.2013.12.030

32.  Fagre AC, Cohen LE, Eskew EA, Farrell M, Glennon E, Joseph MB, et al. Assessing the risk of human-to-wildlife pathogen transmission for conservation and public health. Ecology Letters. 2022;25: 1534–1549. doi:10.1111/ele.14003

33.  Monkeypox – Cameroon. In: World Health Organization: Disease Outbreak News [Internet]. 5 Jun 2018 [cited 13 May 2023]. Available: https://www.who.int/emergencies/disease-outbreak-news/item/05-june-2018-monkeypox-cameroon-en

34.  de Oliveira JS, Figueiredo P de O, Costa GB, de Assis FL, Drumond BP, da Fonseca FG, et al. Vaccinia Virus Natural Infections in Brazil: The Good, the Bad, and the Ugly. Viruses. 2017;9: 340. doi:10.3390/v9110340

35.  Figuerola J, Green AJ. Haematozoan Parasites and Migratory Behaviour in Waterfowl. Evolutionary Ecology. 2000;14: 143–153. doi:10.1023/A:1011009419264

36.  Senkevich TG, Yutin N, Wolf YI, Koonin EV, Moss B. Ancient Gene Capture and Recent Gene Loss Shape the Evolution of Orthopoxvirus-Host Interaction Genes. mBio. 12: e01495-21. doi:10.1128/mBio.01495-21

37.  Carlson CJ, Gibb RJ, Albery GF, Brierley L, Connor RP, Dallas TA, et al. The Global Virome in One Network (VIRION): an atlas of vertebrate-virus associations. 2021; 2021.08.06.455442. doi:10.1101/2021.08.06.455442

38.  Upham NS, Esselstyn JA, Jetz W. Inferring the mammal tree: Species-level sets of phylogenies for questions in ecology, evolution, and conservation. PLOS Biology. 2019;17: e3000494. doi:10.1371/journal.pbio.3000494

39.  Soria CD, Pacifici M, Di Marco M, Stephen SM, Rondinini C. COMBINE: a coalesced mammal database of intrinsic and extrinsic traits. Ecology. 2021;102: e03344. doi:10.1002/ecy.3344

830  40.  The IUCN Red List of Threatened Species. In: IUCN Red List of Threatened Species
831       [Internet]. [cited 14 Jul 2023]. Available: https://www.iucnredlist.org/en

832  41.  Faurby S, Davis M, Pedersen RØ, Schowanek SD, Antonelli1 A, Svenning J-C.
833       PHYLACINE 1.2: The Phylogenetic Atlas of Mammal Macroecology. Ecology. 2018;99:
834       2626–2626. doi:10.1002/ecy.2443

835  42.  Mull N, Carlson CJ, Forbes KM, Becker DJ. Virus isolation data improve host predictions
836       for New World rodent orthohantaviruses. Journal of Animal Ecology. 2022;91: 1290–1302.
837       doi:10.1111/1365-2656.13694

838  43.  Elith J, Leathwick JR, Hastie T. A working guide to boosted regression trees. Journal of
839       Animal Ecology. 2008;77: 802–813. doi:10.1111/j.1365-2656.2008.01390.x

840  44.  De'ath G. Boosted Trees for Ecological Modeling and Prediction. Ecology. 2007;88: 243–
841       251. doi:10.1890/0012-9658(2007)88[243:BTFEMA]2.0.CO;2

842  45.  Benjamini Y, Hochberg Y. Controlling the False Discovery Rate: A Practical and Powerful
843       Approach to Multiple Testing. Journal of the Royal Statistical Society: Series B
844       (Methodological). 1995;57: 289–300. doi:10.1111/j.2517-6161.1995.tb02031.x

845  46.  Yekutieli D, Benjamini Y. Resampling-based false discovery rate controlling multiple test
846       procedures for correlated test statistics. Journal of Statistical Planning and Inference.
847       1999;82: 171–196. doi:10.1016/S0378-3758(99)00041-5

848  47.  Washburne AD, Silverman JD, Morton JT, Becker DJ, Crowley D, Mukherjee S, et al.
849       Phylofactorization: a graph partitioning algorithm to identify phylogenetic scales of
850       ecological data. Ecological Monographs. 2019;89: e01353. doi:10.1002/ecm.1353

851  48.  Holm S. A Simple Sequentially Rejective Multiple Test Procedure. Scandinavian Journal of
852       Statistics. 1979;6: 65–70.

853

854

# Supporting information

856  **S1 Fig. Predictions of *Orthopoxvirus* positivity for the susceptible host model trained on**

857  **virus isolation data.** Bars/segments are scaled by probabilities and colored yellow according to

858  binary classification of predicted host genera assuming (a) an 80% sensitivity threshold versus

859  (b) a threshold maximizing the sum of sensitivity and specificity. Clades identified through

860  phylogenetic factorization with significantly different mean predictions are shaded in grey.

861

862  **S2 Fig. Trait profiles of *Orthopoxvirus* positive mammal genera for the host prediction**

863  **models**. Partial dependence plots for the ten most significant predictors across BRTs (applied

864  to 100 random splits of training and test data) are displayed in order of relative importance for

865  (a) the host exposure model trained on PCR data, and (b) the susceptible host model trained on

866  virus isolation data as the response variable. Grey lines or points represent the marginal effect

867  of each variable on predicting host status for each data split, while black lines indicate the

868  average marginal effect. Histograms and rug plots illustrate the distribution of continuous and

869  categorical predictors, respectively, across all included mammal genera.

870  **S3 Fig. Trait profiles of mammal genera with predicted existence of an *Orthopoxvirus* link**

871  **for the link prediction model.** Partial dependence plots for the top ten predictors across BRTs

872  (applied to 100 random partitions of training and test data) are displayed in order of relative

873  importance for the link prediction model trained on (a) host and viral features and (b) host traits

874  only. Grey lines or points represent the marginal effect of each variable on predicting host-virus

875  links for each data split, while black lines indicate the average marginal effect. Histograms and

876  rug plots display the distribution of continuous and categorical predictors, respectively, across

877  all included mammal genera.

878  **S4 Fig. Marginal effects of the host traits with highest relative importance in the link**

879  **prediction model when trained on host and viral features.** Partial dependence plots for

880  island dwelling (A) and dispersal (B) across BRTs applied to 100 random partitions of training

881  and test data are displayed in order of relative importance.

882  **S5 Fig. Cumulative variance plot for the first ten principal components (PC).** The dashed

883  line displays the cutoff where the principal components capture approximately 70% of the

884  variance in the viral accessory gene data.

885  **S6 Fig. Relative influence of model features ranked for the link prediction model trained**

886  **only on host traits.** Each horizontal bar plots the variability in the relative influence of the

887  predictor variables as measured across 100 random partitions of training (70%) and test (30%)

26

888     data. Boxplots show the median and interquartile range, whiskers are the extremes, and circles

889     are additional outliers.

890     **S7 Fig. Plot of *Orthopoxvirus* accessory gene loadings by the predicted gene function for**

891     **each principal component (PC).** Loading plots with variables (i.e., genes) grouped by their

892     predicted gene function are displayed for (a) PC1 and PC2, (b) PC3 and PC4, (c) PC5 and PC6,

893     (d) PC7 and PC8, and (e) PC9 and PC10. Only genes with loading values greater than or less

894     than 1.5 times the standard deviation from the mean are represented. Genes predicted to play a

895     role in host interaction and immune evasion include those predicted to encode the Ankyrin

896     family of proteins, Kelch-like proteins, and proteins involved in chemokine/cytokine regulation,

897     cell death/cell cycle regulation, pathogen recognition, and antagonizing host adaptive immunity.

898     **S8 Fig. Plot of *Orthopoxvirus* sequences by viral species for each principal component.**

899     Score plots with samples (i.e., *Orthopoxvirus* genome) grouped by viral species are displayed

900     for (a) PC1 and PC2, (b) PC3 and PC4, (c) PC5 and PC6, (d) PC7 and PC8, and (e) PC9 and

901     PC10.

902     **S9 Fig. Geographic distribution of mpox virus hosts.** Host distributions are based on the

903     IUCN Red List database of mammal geographic ranges for those species belonging to (a)

904     observed host genera and (b) predicted (both observed and unobserved) host genera based on

905     the results of the link prediction model and applying a 90% sensitivity threshold for host

906     classification. The corresponding legend depicts the number of species with overlapping

907     geographic range by color.

908     **S10 Fig. Host trait coverage across mammal genera.** Features with at least 60% coverage

909     (denoted by the dashed line) across mammal genera were included in the BRT models. A

910     complete list of feature coverage is available on the GitHub repository.

911     **S11 Fig. Performance measures of boosted regression tree models during parameter**

912     **tuning.** Performance was evaluated by the area under the receiver operating characteristic

913     curve (AUC), sensitivity, and specificity for (a) the host exposure model trained on RT-PCR

914  data, (b) the susceptible host model trained on virus isolation data, and (c) the link prediction

915  model during parameter tuning. Boxplots show the median and interquartile range alongside

916  raw data for all 10 random splits of training (70%) and test (30%) data for each combination of

917  learning rate and interaction depth.

918  **S1 Table. Mean performance measures for each ensemble model and their**

919  **corresponding standard errors (SE).** Performance was evaluated by the area under the

920  receiver operating characteristic curve (AUC), the sensitivity, and the specificity for each model

921  assuming a threshold value of 0.5. We used 100 random partitions to generate an ensemble. An

922  unpaired *t*-test with *p*-values adjusted for the false discovery rate using the Benjamini Hochberg

923  correction compares the performance of the host exposure model (based on molecular

924  detection of viral DNA using PCR techniques) to the susceptible host model (based on virus

925  isolation from hosts) and that of the link prediction models with vs. without vaccinia virus

926  associations. The effect size for the *t*-test was calculated using Cohen's *d,* which standardizes

927  the mean difference.

928  **S2 Table. Phylogenetic factorization of mean predicted probabilities for *Orthopoxvirus***

929  **positivity for the host exposure model.** The number of retained clades after a 5% family-wise

930  error rate, taxa corresponding to those clades, number of species per clade, and mean

931  predicted probabilities for the clade compared to the paraphyletic remainder are shown.

932  **S3 Table. Phylogenetic factorization of mean predicted probabilities for *Orthopoxvirus***

933  **positivity for the susceptible host model.** The number of retained clades after a 5% family-

934  wise error rate, taxa corresponding to those clades, number of species per clade, and mean

935  predicted probabilities for the clade compared to the paraphyletic remainder are shown.

936  **S4 Table. Phylogenetic factorization of mean predicted probabilities for host-**

937  ***Orthopoxvirus* associations for the link prediction model.** The number of retained clades

938  after a 5% family-wise error rate, taxa corresponding to those clades, number of species per

939   clade, and mean predicted probabilities for the clade compared to the paraphyletic remainder

940   are shown.

941   **S5 Table. Estimation of threshold values ($t$) based on different optimal thresholding**

942   **methods.** Changes in the number of predicted hosts ($n$), sensitivity, and specificity were

943   calculated for each threshold value applied to each ensemble model.

944   **S6 Table. Importance and ranks of mammal traits for the host prediction models.** The host

945   exposure model (based on molecular detection of viral genomes using PCR techniques) and the

946   susceptible host model (based on virus isolation from hosts) were trained on 105 predictor

947   variables, including but not limited to 61 ecological features and 42 taxonomic trait variables.

948   **S7 Table. Importance and rank of mammal traits and viral features for the link prediction**

949   **model.** The link prediction model was trained on 183 predictor variables, including but not

950   limited to 10 viral features (i.e., principal components), 61 ecological traits, and 110 taxonomic

951   traits.

952   **S1 Data. Predictive function of associated proteins for genes ranked by their positive and**

953   **negative loading contribution to principal components 1, 4, 3 and 9.** Genes with loading

954   values that were greater than the mean plus 1.5 times the standard deviation and less than the

955   mean minus 1.5 times the standard deviation were included for those with positive loadings

956   separate from those with negative loadings.

957   **S2 Data. Host-virus associations for link prediction based on *Orthopoxvirus* genomes**

958   **extracted from NCBI.**

959