

1 **Machine learning approaches for estimating cross-neutralization potential among FMD**  
2 **serotype O viruses**

3 **Running title: *In silico* prediction of r<sub>1</sub> values in FMDVs**

4 Dennis N Makau <sup>a, c#</sup>, Jonathan Arzt <sup>b</sup>, Kimberly VanderWaal <sup>c</sup>

5 <sup>a</sup> Department of Biomedical and Diagnostic Sciences, College of Veterinary Medicine, University  
6 of Tennessee, Knoxville, USA.

7 <sup>b</sup> Foreign Animal Disease Research Unit, USDA-ARS, Plum Island Animal Disease Center, Orient  
8 Pt., NY 11957, USA

9 <sup>c</sup> Department of Veterinary Population Medicine, College of Veterinary Medicine, University of  
10 Minnesota, USA

11

12

13 #Address correspondence to Dennis N Makau, [dmakau@utk.edu](mailto:dmakau@utk.edu)

14

15

16 **Abstract**

17 In this study, we aimed to develop an algorithm that uses sequence data to estimate cross-  
18 neutralization between serotype O foot-and-mouth disease viruses (FMDV) based on r<sub>1</sub> values,  
19 while identifying key genomic sites associated with high or low r<sub>1</sub> values. The ability to estimate  
20 cross-neutralization potential among co-circulating FMDVs in silico is significant for vaccine  
21 developers, animal health agencies making herd immunization decisions, and disease  
22 preparedness. Using published data on virus neutralization titer (VNT) assays and associated VP1  
23 sequences from GenBank, we applied machine learning algorithms (BORUTA and random forest)  
24 to predict potential cross-reaction between serum/vaccine-virus pairs for 73 distinct serotype O  
25 FMDV strains. Model optimization involved tenfold cross-validation and sub-sampling to address  
26 data imbalance and improve performance. Model predictors included amino acid distances, site-  
27 wise amino acid polymorphisms, and differences in potential N-glycosylation sites.

28 The dataset comprised 108 observations (serum-virus pairs) from 73 distinct viruses with r<sub>1</sub>  
29 values. Observations were dichotomized using a 0.3 threshold, yielding putative non-cross-  
30 neutralizing (< 0.3 r<sub>1</sub> values) and cross-neutralizing groups (≥ 0.3 r<sub>1</sub> values). The best model had  
31 a training accuracy, sensitivity, and specificity of 0.96 (95% CI: 0.88-0.99), 0.93, and 0.96,  
32 respectively, and an accuracy of 0.94 (95% CI: 0.71-1.00), sensitivity of 1.00, and specificity of  
33 0.93, positive, and negative predictive values of 0.60 and 1.00, respectively, on one testing dataset  
34 and an accuracy, AUC, sensitivity, specificity, and predictive values all approaching 1.00 on a  
35 second testing dataset. Additionally, amino acid positions 48, 100, 135, 150, and 151 in the VP1  
36 region alongside amino acid distance were found to be important predictors of cross-neutralization.

37 Our study highlights the value of genetic/genomic data for informing immunization strategies in  
38 disease management and understanding potential immune-mediated competition amongst related

39 endemic strains of serotype O FMDVs in the field. We also showcase leveraging routinely  
40 generated sequence data and applying a parsimonious machine learning model to expedite  
41 decision-making in selection of vaccine candidates and application of vaccines for controlling  
42 FMD, particularly serotype O. A similar approach can be applied to other serotypes.

43 Key words: cross-protection;  $r_1$  values, cross-reactivity; viral neutralization; machine learning;  
44 bioinformatics; computational immunology

## 45 **Introduction**

46 Foot and mouth disease (FMD) is a viral disease caused by the FMD virus (FMDV), a member of  
47 the *Picornaviridae* family (1) that affects cloven-hoofed ungulates. Though typically not fatal, the  
48 impact of FMD on food security and livelihoods in endemic countries (particularly low-and  
49 middle-income countries (LMICs)) cannot be overemphasized (2). FMD also continues to be a  
50 major stumbling block to livestock production and global trade in many parts of the world,  
51 especially with the continued threat of emergence and introduction of new FMDV lineages/strains  
52 into FMD-free countries. As such, numerous efforts and pathways to achieve and maintain disease  
53 free status in most countries has required, among other things, the ability and capacity to identify  
54 viruses and vaccinate animals with appropriate vaccines for countries where elimination is yet to  
55 be achieved. Although there are 7 documented antigenically distinct serotypes (O, A, C, SAT1,  
56 SAT2, SAT3 and Asia1), serotype A and O have been reported to be the most common causes of  
57 FMD globally (3–6).

58 Continued efforts for disease management have included the identification of suitable vaccine  
59 candidates and formulation of vaccines applied to the animal populations at risk. To identify these  
60 vaccine candidates, a process known as vaccine matching is done, which entails the use of viral  
61 neutralization titer (VNT) assays to identify effective cross-reaction and cross-neutralization of  
62 viruses by candidate vaccines. The measure of effectiveness of the cross-neutralization and  
63 subsequent cross-protection is expressed as the  $r_1$  value and is briefly defined as the ratio between  
64 the neutralizing serum titer against the heterologous virus (usually a field strain) to the neutralizing  
65 serum titer against the homologous virus (usually the vaccine strain). Subsequently,  $r_1$  values  
66 between virus-serum pairs of  $\geq 0.3$  are considered evidence of cross-neutralization between the  
67 two involved viruses, while  $r_1$  values  $< 0.3$  are considered to be non-neutralizing (7).

68 Generating  $r_1$  data for comparison demands significant effort, precision, and laboratory resources.  
69 This involves producing sera in live animals and running *in vitro* assays, to evaluate the cross-  
70 reactivity of target viruses with various sera samples. Advances in computational biology empower  
71 us to harness machine learning and artificial intelligence to build computation tools aimed at  
72 estimating the cross-neutralization potential between different viruses. While previous studies  
73 have explored predictive models of FMDV antigenicity, they have varied objectives. Some  
74 focused on predicting FMDV antigenicity and identifying antigenicity descriptors (8), identifying  
75 and predicting specific epitopes and genomic regions of antigenic importance (9, 10), or  
76 developing intricate models to forecast antigenic changes among SAT 1 and 2 FMDV serotypes  
77 (11–13). There is hardly any literature specific to serotype O despite it being one of two widely  
78 occurring serotypes, nor has research on the application of machine learning to estimate  $r_1$  values  
79 been published. Although there have been some divergent views about the adequacy of using  $r_1$

80 values in identifying potential cross-protection (14), it is still one of the metrics relied upon when  
81 interpreting data generated from *in vitro* VNT assays and selecting candidate vaccine strains.  
82 Successful application of machine learning to aid in decision making when selecting vaccine  
83 candidates and immunization programs has been demonstrated in diseases such as influenza and  
84 dengue virus (15–17). Our intention in this study was to develop a simple, yet robust, predictive  
85 tool that can be used to estimate potential cross-neutralization between viruses by leveraging  
86 machine learning algorithms and genetic characteristics of FMDV.

87 Therefore, the objective of this study was to develop an algorithm to estimate potential cross-  
88 neutralization between strains of serotype O FMDV using  $r_1$  values and identify important genomic  
89 sites that influence high or low  $r_1$  values between viruses. The ability to distinguish and estimate  
90 the potential for cross-protection between different cocirculating FMD viruses *in silico* will  
91 support decision-making by vaccine developers in vaccine candidate matching, animal health  
92 agencies/institutions in decision-making for herd immunization practices, and overall disease  
93 preparedness and response in cases of emergent strains especially for serotype O about which there  
94 is a dearth of information.

## 95 **Materials and methods**

### 96 **Data preparation**

97 We obtained data from published manuscripts on  $r_1$  values for FMDV serotype O. Viruses utilized  
98 in these papers came from 14 countries. These manuscripts were obtained from a search from  
99 google and PubMed databases using the search terms ‘FMDV and  $r_1$  values’ and ‘FMDV and  
100 vaccine matching’ yielding 4 manuscripts (18–21). To be included, the study must have followed  
101 standard WHOA (22) methodology for generation of  $r_1$  values for serotype O viruses using virus  
102 neutralization tests (VNTs), as well as have associated viral sequence data. The  $r_1$  values indicate  
103 the serological compatibility between the vaccine strain and the field isolate, determined by  
104 comparing the reactivity of the heterologous versus homologous virus to the antisera produced  
105 against the specific viruses. The ratio of the heterologous to homologous neutralization titers  
106 between field isolates and vaccine sera serve as a measure of cross-protection.

107 Since majority of the published studies had uploaded only VP1 data to GenBank, we downloaded  
108 sequence data for the VP1 gene of 73 distinct viruses used in the cross-reactivity and generation  
109 of  $r_1$  data summarized in the manuscripts described above and any whole genomes available were  
110 trimmed to the VP1 region using Muscle in Aliview 1.26 (23). Accession numbers to the sequences  
111 used in this analysis are available in supplemental material 1. Using Muscle in Aliview 1.26 (23)  
112 we aligned and translated sequences into amino acid sequences. Using MEGA X (24), we  
113 calculated Poisson corrected amino acid distance between all pairs (24, 25) and used R packages  
114 *stringr* (26), *bioseq* (27), *ape* (28) and *tidysq* (29) to code site-wise differences at any polymorphic  
115 amino acid sites using R v4.2 software (30). For each polymorphic amino acid site, serum-virus  
116 pairs were tabulated as 0 if they shared the same amino acid and 1 if they had a different amino  
117 acid at that site. We also scanned the genome for sites of potential N-glycosylation (31, 32) in the  
118 VP1 region using a custom-built R script. Pairs with identical sets of inferred glycosylated sites  
119 were coded as 0 and non-identical sets as 1. The  $r_1$  values, as reported in the manuscripts, amino  
120 acid distance calculated in MEGA X, and amino acid site-wise differences were concatenated into  
121 one data frame, resulting in 108 observations (serum-virus pairs). Although some studies on  
122 influenza have explored the pros and cons of weighting of specific amino acid sites based on a

123 *priori* knowledge of their evolutionary history and influence on antigenicity (33, 34), our modeling  
124 approach was naïve to *a priori* assumptions about the relative importance of specific site-wise  
125 amino acid differences and thus no weighting was used (e.g.,  $I \rightarrow N \equiv N \rightarrow I$ ) (15).

126 Initial analysis involved descriptive statistics of the data and correlation analysis. We performed a  
127 phylogenetic analysis of the sequence data by constructing a bootstrapped maximum likelihood  
128 tree in RAxML with 500 bootstraps to depict the distribution/representation of the different  
129 topotypes in our data and to visualize relatedness of the serum/vaccine candidates cross-reacted in  
130 VNT assays to generate  $r_1$  values. Additionally, we performed a Spearman's correlation analysis  
131 between p-amino acid distance and  $r_1$  scores. The  $r_1$  values were then dichotomized at 0.3, which  
132 is the threshold recommended by WHOA (22).

### 133 **Training and testing machine learning models**

134 Using a stepwise approach, we developed a random forest classification model with sub-sampling  
135 steps and tenfold cross validation to achieve the best model performance. Upon dichotomization  
136 of the  $r_1$  values into binomial data, there was an imbalance in the 0 (non-cross-reacting) vs 1 (cross-  
137 reacting) classes with more than twice as many observations in the 0 category as those in the 1  
138 category (0 = 87 observations, 1 = 21 observations). As such, there was need to subsample the data  
139 to adjust for this imbalance of the outcome in the training dataset; we thus used the synthetic  
140 minority oversampling technique (SMOTE) which statistically increased the number of cases in  
141 the dataset to balance the distribution of cases vs non-cases (1 vs 0) in the outcome to achieve  
142 better model performance. Feature-reduction was then implemented to reduce the number of model  
143 features (i.e., predictors) using the Boruta algorithm to optimize model fitting. Upon concatenating  
144 all parts of the data into a single data frame, the final data frame had 216 columns (amino acid  
145 distance and 214 site-wise amino acid differences and potential sites for N-glycosylation). With  
146 highly dimensional data, machine learning models may struggle to achieve good model accuracy  
147 and prediction due to noise introduced as the model tries to optimize the estimated contribution of  
148 each model feature to the variation in the outcome ( $r_1$  class in our study). As such, the Boruta  
149 algorithm has been proposed as a way of eliminating correlated and redundant features from the  
150 model, thus reducing the number of model features needed in the final model and optimizing model  
151 performance. After parsing the data through Boruta, we included all model features classified as  
152 important or tentatively important, resulting in a data frame with 35 model features. These features  
153 were VP1 amino acid distance and site-wise differences in amino acid positions (4, 13, 24, 32, 33,  
154 43, 48, 49, 57, 69, 97, 100, 124, 135, 139-141, 143, 145, 149-151, 154, 156, 159, 166, 173, 175,  
155 195, 198, 199, 210, 213, 214) and potential N-glycosylation profile.

156 Subsequently we used the 80:20 split rule to randomly split the data into three subsets to generate  
157 one training and two test datasets with 70, 17 and 21 observations each. Using *caret* (35) and  
158 *randomForest* (36) packages in R for model building, training and evaluation, including hyper-  
159 parameter tuning and 10-fold cross-validation, we trained the model on 500 trees, tested it, and  
160 made predictions as summarized in the results section. We based model performance on accuracy  
161 (overall percent of observations correctly classified as cross-reacting/non-cross-reacting),  
162 sensitivity (percent of high observations correctly classified), specificity (percent of low  
163 observations correctly classified), and predictive values (proportion (%) of times the classification  
164 (non-cross-protecting-0 vs cross-protecting-1) is the true  $r_1$  group class.

165 Additionally, from the random forest model, model features were ranked in their importance to the  
166 performance of our model using mean decrease in model accuracy when a features' data were

167 randomized relative to the outcome (the relative prediction strength of a variable) and  
168 improvement in Gini index (measure of node impurity associated with a variable) when data were  
169 split on a variable (37, 38). This allowed us to highlight highly ranked amino acids in our VP1 data,  
170 and considered the relative importance of different amino acid sites based on their role in  
171 improving model accuracy and node purity in outcome classification (15, 39).

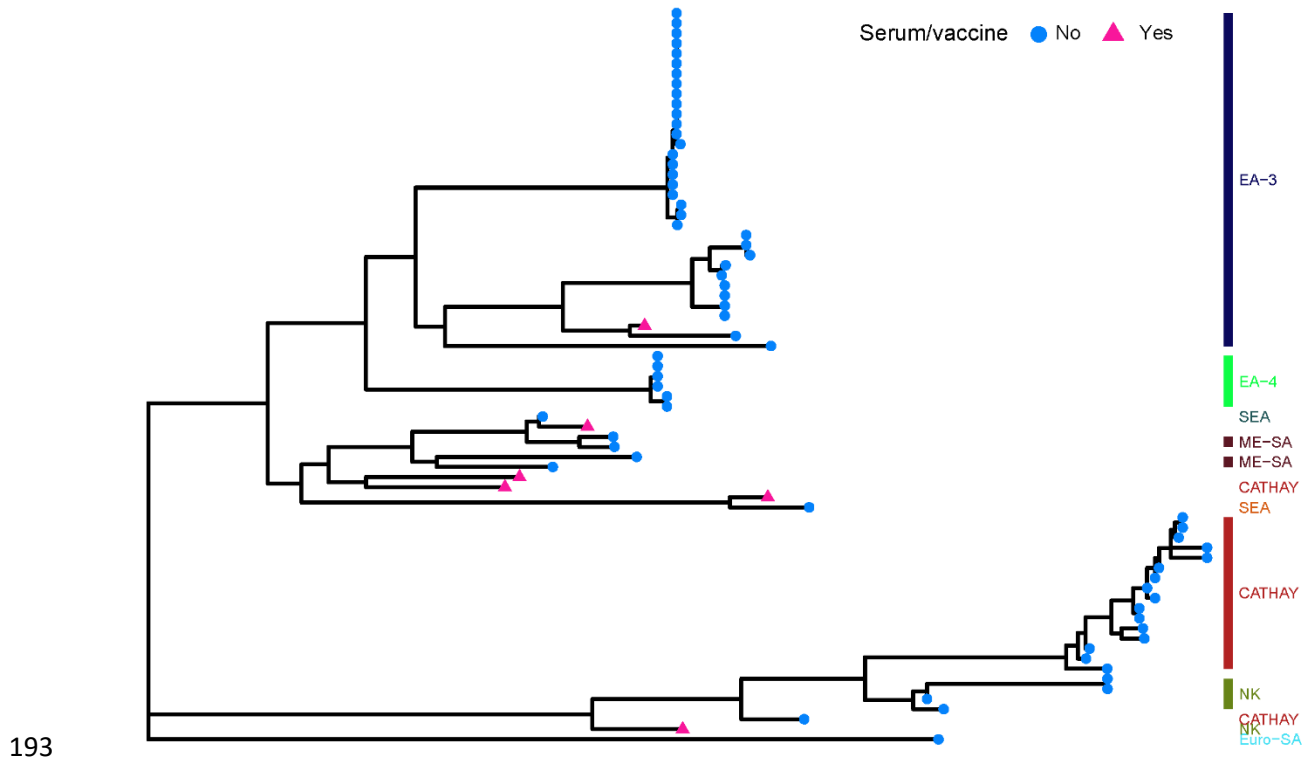
172 To test the performance of the machine learning algorithm's predictions on a completely external  
173 set of sequences that have been reported to be antigenically novel, we used data from Bachanek-  
174 Bankowska et al., (40). In this study, they reported the isolation of three antigenically distinct  
175 viruses isolated from outbreaks in Pakistan in 2016 and 2017. Bachanek-Bankowska et al.,  
176 described three isolates (GenBank accession numbers MH784403, MH784404, and MH784405)  
177 from a single genetic sublineage displaying distinct antigenic phenotypes against three commonly  
178 used vaccines in the region (O 3039 and O Manisa [Boehringer Ingelheim] and O TUR 5/2009  
179 [MSD]).

## 180 **Results**

### 181 *Descriptive analysis*

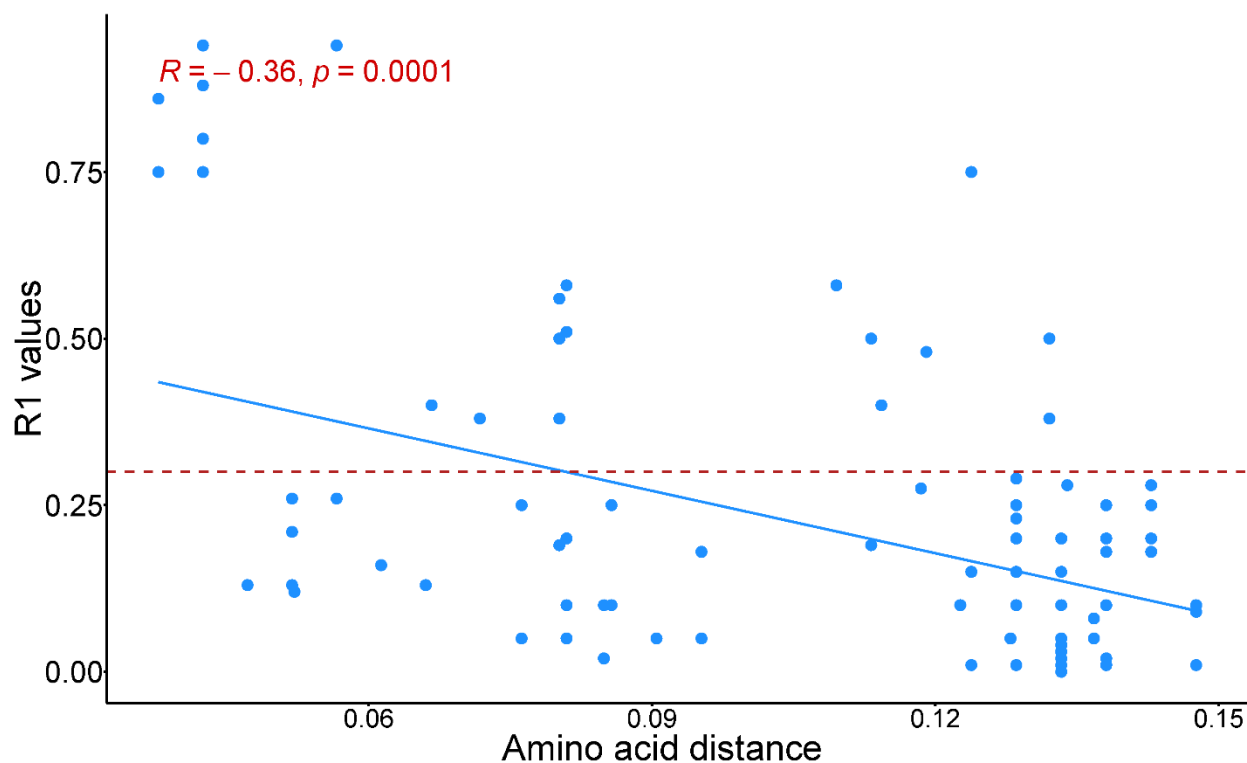
182 In this dataset, there were 108 observations (serum-virus pairs) from 73 distinct viruses for which  
183  $r_1$  values had been reported. In summary, the mean  $r_1$  value was  $0.22 \pm 0.23$  (SD), with a range of  
184 (0.0-0.94) and a median of 0.16. Upon dichotomizing the data using the 0.3 threshold as cutoff,  
185 the two resulting groups used in subsequent modelling comprised of 87 pairs in the non-cross-  
186 neutralizing group ( $< 0.3$   $r_1$  values) and 21 in the cross-neutralizing group ( $\geq 0.3$   $r_1$  values).

187 Additionally, upon phylogenetic assessment, our dataset included viruses belonging to five  
188 topotypes for FMDV serotype O and the serum/vaccines used in the assays were evenly distributed  
189 among the topotype groupings (Figure 1). The mean VP1 p-amino acid distance for the pairs was  
190  $0.12 \pm 0.03$  with a range of 0.04 – 0.15 and median of 0.13. There was a significant ( $p=0.0001$ )  
191 moderate negative correlation (Spearman  $\rho = -0.4$ ) between amino acid distance and  $r_1$  values  
192 between the pairs (Figure 2).



194 Figure 1: Maximum likelihood tree depicting the phylogenetic distribution of 73 foot and mouth  
195 viruses (serum, vaccine (pink triangles) and field strains (blue circles)) representing serotype O  
196 topotypes used in the development of an  $r_1$  predictive model. The colored bars on the right of the  
197 tree represent the different topotypes included in the data (*NK=Not Known*).

198



199

200 Figure 2: Correlation plot depicting an inverse moderate Spearman's correlation (blue) between  
201 amino acid distance of the VP1 region for serotype O foot and mouth disease viruses and  $r_1$  values  
202 obtained from *in vitro* virus neutralization assays. The red-dotted line represents the cut-off of 0.3  
203 threshold for potentials cross-neutralization between serum-virus pairs.

#### 204 **Random forest model performance and predictions**

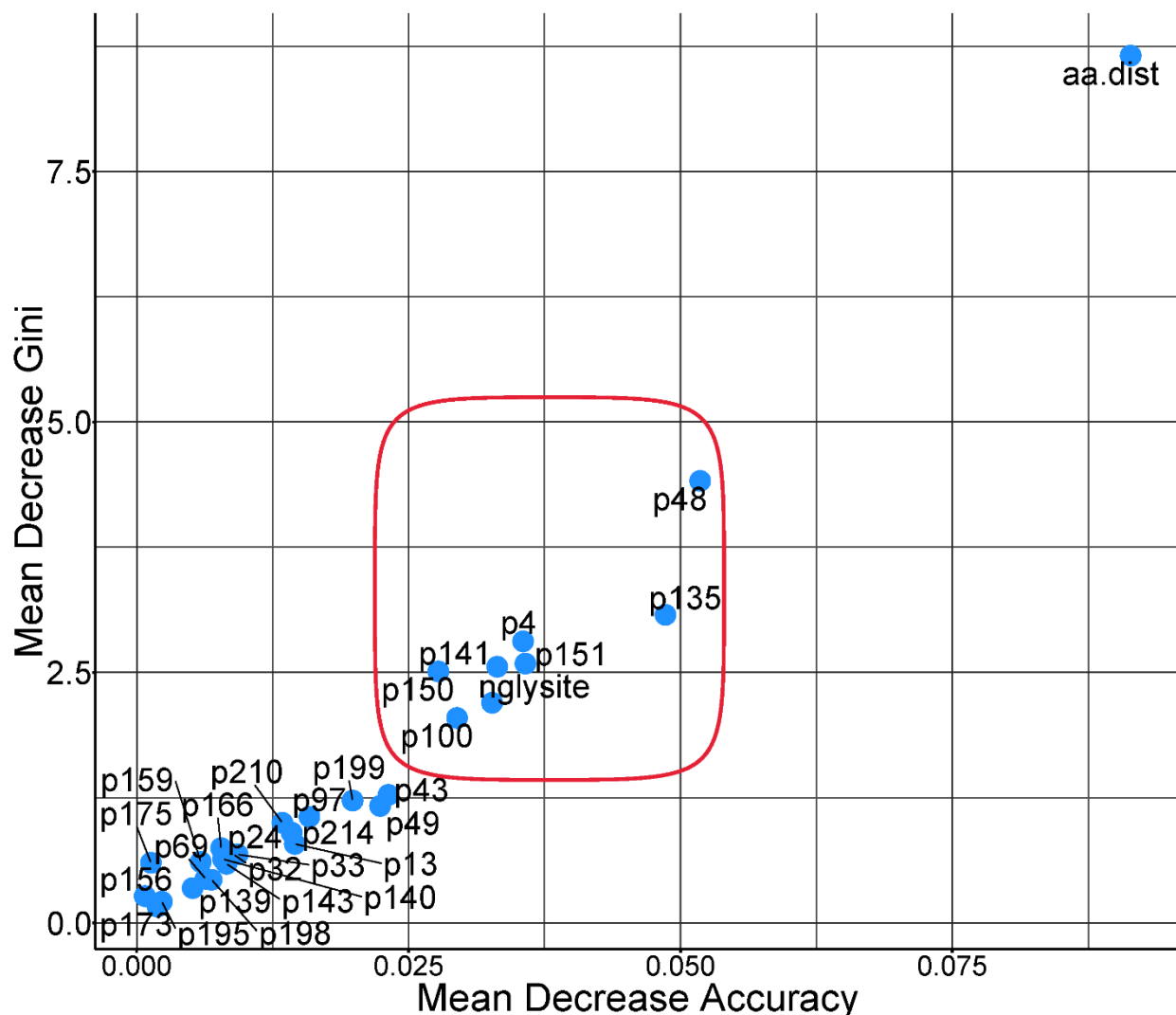
205 The best performing model had an accuracy of 0.96 (95% CI 0.88-0.99), AUC of 0.95, sensitivity  
206 and specificity of 0.93 and 0.96, and negative and positive predictive values of 0.98 and 0.87,  
207 respectively on training data (n=70: 0=56, 1=14). On one testing dataset (n=17: 0=14, 1=3), the  
208 model accuracy was 0.88 (95% CI 0.64-0.99), AUC of 0.93, sensitivity of 1.00 and specificity of  
209 0.86, and positive and negative predictive values of 0.60 and 1.00, respectively. The model  
210 performance on the second test/validation data (n=21: 0=17, 1=4) was 1.00 accuracy (95% CI  
211 0.84-1.00), AUC of 1.00, sensitivity and specificity of 1.00 and 1.00, and positive and negative  
212 predictive values of 1.00 and 1.00, respectively.

213 From the variable importance analysis that extracted which predictors had the most influence in  
214 the models' predictive ability and accuracy, several sites were ranked highly besides the amino  
215 acid distance. These include amino acid positions 48, 100, 135, 150, and 151 in the VP1 region  
216 (Figure 3).

217 When the trained model was applied to the antigenically novel sequences described by Bachanek-  
218 Bankowska (40), the model correctly predicted the antigenic relationship of one out of three  
219 antigenically distinct viruses, i.e., between MH784405.1 (PAK/14/2017) and commonly used  
220 vaccines in the region. The model however incorrectly predicted that commonly used vaccines in

221 the regions were good matches for both MH784403.1 (PAK/10/2016) and MH784404.1  
222 (PAK/4/2017) ( $r_1 \geq 0.3$ ). The *in vitro* vaccine matching experiment had indicated non-cross-  
223 neutralization between the two outbreak viruses and commonly used vaccines in this region  
224 (Supplementary table 1).

225



226  
227 Figure 3: Multiway importance plot for model covariates included in the final random forest model  
228 predicting  $r_1$  values in VP1 region for serotype O foot and mouth disease virus. Feature labels in  
229 the graph indicate the specific amino acid site (p) while 'aa.dist' indicates the amino acid distance  
230 between serum-virus pairs, 'nglysite' similarity or difference in positions for potential N-  
231 glycosylation.

232

233 **Discussion**



234 We developed a parsimonious random forest classification model that estimates potential cross-  
235 neutralization between serum and viruses i.e.,  $r_1$  values for viruses belonging to serotype O FMDV  
236 with an accuracy of more than 85%. This model adds to a growing body of research and efforts  
237 towards leveraging bioinformatic data to streamline and enhance our understanding on viral  
238 antigenicity and host immune response as exemplified in diseases like influenza (16, 17, 41–44).  
239 diseases like influenza. Optimized and robust predictive models can be valuable in decision  
240 making and implementing strategic interventions towards the control, and ultimately eradication,  
241 of FMDV. This model requires further training and refinement with diverse datasets because when  
242 applied to the antigenically novel sequences described by Bachanek-Bankowska (40), the model  
243 correctly predicted the antigenic relationship of only one of three novel viruses.

244 The suitability of using *in vitro* methods to estimate *in vivo* cross-protection in FMD (45),  $r_1$  values  
245 and 0.3 thresholds compared to other measures of immune response and cross reactivity has been  
246 argued by several studies (14, 45–48). However, though this model is based on  $r_1$  values it can be  
247 modified to accommodate other measures of antigenic variability like raw VNT results. Since  $r_1$   
248 values are still the most-used option when conducting vaccine matching and selection of  
249 candidates, whether achieved through simple ELISA or a modified combination of techniques to  
250 estimate cross-protection, our model is well trained to aid in the identification of potential cross-  
251 neutralizing and non-cross-neutralizing sera thus streamlining the selection of potential vaccine  
252 candidates and facilitating immediate comparisons and of field strains where necessary. Moreover,  
253 when applied in non-research settings, animal health agencies/institutions can easily upload  
254 sequences to the model from isolates or clinical samples and promptly identify potential options  
255 for immunization. The continued improvement and accessibility of sequencing technology such as  
256 Minion Nanopore sequencing, which can be deployed in LMIC settings, may also increase the  
257 frequency and capacity of in-country generation of sequences (especially in countries where access  
258 to reference labs is limited) complementing the utility of predictive tools supported by this model.

259 By highlighting the highly ranked/important amino acids in the VP1 region, we posit that those  
260 amino acids influence host immune induction and response as differences in those sites appear to  
261 be important for estimating the  $r_1$  values between serum-virus pairs. This supports findings by  
262 other studies (10, 12) that certain immunodominant sites exist in the VP1 region. To better  
263 understand challenges with vaccine effectiveness, causes of vaccine failure and breakthrough  
264 outbreaks in vaccinated populations as well as the mechanisms of emergence and spread of new  
265 viruses, dynamics of mutations in these sites can be investigated further (49). In these data, seven  
266 amino acid positions were most influential in the ability of our model to accurately distinguish  
267 between cross-neutralizing ( $\geq 0.3$ ) vs non-cross-neutralizing ( $< 0.3$ ). These were amino acids in  
268 positions 4, 48, 100, 135, 141, 150 and 151. Except position 4 and 100, all these amino acids are  
269 part of immunogenically important regions described as the B-C and G-H loops in the VP1 protein  
270 of serotype O (6) in which certain mutations in combination or independently have been thought  
271 to confer serological heterogeneity among FMDVs (50). In sequences used for this study, N-  
272 glycosylation was predicted to occur most often at sites 86, then 101, and least often at site 132.

273 Even with improved predictive ability to accurately classify cross-neutralizing sera/vaccine as  
274 possible, there are drawbacks of *in silico* modelling for biological processes. The immune response  
275 is a composite process which commonly involves the interplay of multiple compartments of the  
276 host immune response system and different factors (e.g., environment and host genetics) that may  
277 influence host-virus interactions. Such factors contribute to why measures of cross-reactivity *in*  
278 *vitro* often do not translate well to *in vivo* cross-protection in the field. Our model does not account

279 for these factors. Nevertheless, with current data, we were able to accurately categorize 8 to 9 out  
280 of every 10 serum-virus pairs for serotype O as cross-neutralizing or non-cross-neutralizing (based  
281 on  $r_1$  values). However, the wider confidence interval suggests that there is some misclassification  
282 and misidentification of sera/vaccine candidates.

283 As such, continued improvement of *in silico* models is necessary both with diverse datasets and  
284 potential validation with *in vivo*/field application data as attempted in this study. That limitation  
285 notwithstanding, with the limited data used here, the model can at least aid in filtering potential  
286 candidates to be considered for in-depth *in vitro/in vivo* assays or blends in vaccination programs  
287 which would save resources and increase efficiency of the decision-making process and FMD  
288 management.

289 Lastly, other VP regions in the P1 portion of the genome also likely play a role in host-virus  
290 interaction and immune response. However, since we only relied on publicly available secondary  
291 data, we did not have access to full genome sequences for all 73 viruses included in the study  
292 hence the need to restrict the analysis to the role of VP1 in estimating  $r_1$  values. Whole genomic  
293 data could improve our model. Future efforts will explore the benefits of a broader genomic  
294 analysis and more sensitive threshold of cross-protection.

## 295 **Conclusion**

296 Ultimately, this study adds to the growing body of literature that machine learning can be applied  
297 to genetic/genomic data to achieve a more nuanced analysis of the relationship between genetic  
298 variability and cross-recognition of viruses by host immune systems. Such models can support  
299 immunization as a pathway to disease management. In this study, we demonstrate the opportunities  
300 and potential for leveraging routinely generated sequence data and the application of a  
301 parsimonious machine learning model to streamline the process of decision-making in vaccine  
302 development and application to control FMD, especially serotype O (but also applicable to other  
303 serotypes). Specifically, by obtaining an accurate model with high sensitivity and specificity,  
304 appropriate vaccine candidates may be able to be selected more quickly, although *in vivo*  
305 experiments would still ultimately be necessary to assess cross-protection. Also, this capability  
306 would enable tailoring immunization protocols for use in the field with a faster turnaround time  
307 for decision-making. Lastly, outputs from such a model can be combined with other mathematical  
308 models to understand drivers and trends of viral emergence, especially the role of immune pressure  
309 in driving the evolution and spread of FMDV. The latest version of the  $r_1$  predictive model is  
310 available for access via a Shiny dashboard (<https://dmakau.shinyapps.io/PredImmune-FMD/>).

311  
312 **Data availability.** The nucleotide sequences of the FMDV used in this analysis are already  
313 publicly available on GenBank and the list provided in Supplemental material 1 can help in  
314 identification and downloading of the sequences if needed.

## 315 **Supplemental material**

316 All supplementary materials have been provided in one file.

## 317 **Acknowledgements**

318 This project was supported by the USDA Agricultural Research Service, award 58-8064-2-006.

319 The authors declare there are no conflicts of interest.

320

## 321 **References**

- 322 1. Zell R, Delwart E, Gorbalenya AE, Hovi T, King AMQ, Knowles NJ, Lindberg AM,  
323 Pallansch MA, Palmenberg AC, Reuter G, Simmonds P, Skern T, Stanway G, Yamashita  
324 T. 2017. ICTV virus taxonomy profile: Picornaviridae. *J Gen Virol* 98:2421–2422.
- 325 2. Knight-Jones TJD, McLaws M, Rushton J. 2017. Foot-and-Mouth Disease Impact on  
326 Smallholders - What Do We Know, What Don't We Know and How Can We Find Out  
327 More? *Transbound Emerg Dis* 64:1079.
- 328 3. He Y, Li K, Cao Y, Sun Z, Li P, Bao H, Wang S, Zhu G, Bai X, Sun P, Liu X, Yang C,  
329 Liu Z, Lu Z, Rao Z, Lou Z. 2021. Structures of Foot-and-mouth disease virus with  
330 neutralizing antibodies derived from recovered natural host reveal a mechanism for cross-  
331 serotype neutralization. *PLoS Pathog* 17:1–20.
- 332 4. Brito BP, Rodriguez LL, Hammond JM, Pinto J, Perez AM. 2017. Review of the Global  
333 Distribution of Foot-and-Mouth Disease Virus from 2007 to 2014. *Transbound Emerg Dis*  
334 64:316–332.
- 335 5. Grubman MJ, Baxt B. 2004. Foot-and-mouth disease. *Clin Microbiol Rev* 17:465–493.
- 336 6. Ranaweera LT, Wijesundara UK, Jayarathne HSM, Knowles N, Wadsworth J, Mioulet V,  
337 Adikari J, Weebadde C, Sooriyapathirana SS. 2019. Characterization of the FMDV-  
338 serotype-O isolates collected during 1962 and 1997 discloses new topotypes, CEY-1 and  
339 WCSA-1, and six new lineages. *Sci Reports* 2019 9:1–10.
- 340 7. OIE/FAO Foot-and-Mouth Disease Reference Laboratories Network. 2022. Foot-and-  
341 Mouth Disease European July-September 2022 Quarterly report.
- 342 8. Qiu J, Qiu T, Dong Q, Xu D, Wang X, Zhang Q, Pan J, Liu Q. 2021. Predicting the  
343 Antigenic Relationship of Foot-and-Mouth Disease Virus for Vaccine Selection through a  
344 Computational Model. *IEEE/ACM Trans Comput Biol Bioinforma* 18:677–685.
- 345 9. Borley DW, Mahapatra M, Paton DJ, Esnouf RM, Stuart DI, Fry EE. 2013. Evaluation  
346 and Use of In-Silico Structure-Based Epitope Prediction with Foot-and-Mouth Disease  
347 Virus. *PLoS One* 8.
- 348 10. Bari FD, Parida S, Asfor AS, Haydon DT, Reeve R, Paton DJ, Mahapatra M. 2015.  
349 Prediction and characterization of novel epitopes of serotype A foot-and-mouth disease  
350 viruses circulating in East Africa using site-directed mutagenesis. *J Gen Virol* 96:1033–  
351 1041.
- 352 11. Rahman T, Mahapatra M, Laing E, Jin Y. 2015. Evolutionary non-linear modelling for  
353 selecting vaccines against antigenically variable viruses. *Bioinformatics* 31:834–840.
- 354 12. Davies V, Reeve R, Harvey WT, Maree FF, Husmeier D. 2017. A sparse hierarchical  
355 Bayesian model for detecting relevant antigenic sites in virus evolution. *Comput Stat*  
356 32:803–843.
- 357 13. Reeve R, Blignaut B, Esterhuysen JJ, Opperman P, Matthews L, Fry EE, de Beer TAP,  
358 Theron J, Rieder E, Vosloo W, O'Neill HG, Haydon DT, Maree FF. 2010. Sequence-

- 359 Based Prediction for Vaccine Strain Selection and Identification of Antigenic Variability  
360 in Foot-and-Mouth Disease Virus. *PLOS Comput Biol* 6:e1001027.
- 361 14. Brito BP, Perez AM, Capozzo AV. 2014. Accuracy of traditional and novel serology tests  
362 for predicting cross-protection in foot-and-mouth disease vaccinated cattle. *Vaccine*  
363 32:433–436.
- 364 15. Zeller MA, Gauger PC, Arendsee ZW, Souza CK, Vincent AL, Anderson TK. 2021.  
365 Machine Learning Prediction and Experimental Validation of Antigenic Drift in H3  
366 Influenza A Viruses in Swine. *mSphere* 6.
- 367 16. Bell SM, Katzelnick L, Bedford T. 2019. Dengue genetic divergence generates within-  
368 serotype antigenic variation, but serotypes dominate evolutionary dynamics. *Elife* 8.
- 369 17. Mansfield KL, Horton DL, Johnson N, Li L, Barrett ADT, Smith DJ, Galbraith SE,  
370 Solomon T, Fooks AR. 2011. Flavivirus-induced antibody cross-reactivity. *J Gen Virol*  
371 92:2821–2829.
- 372 18. Maree FF, Blignaut B, Esterhuysen JJ, de Beer TAP, Theron J, O’Neill HG, Rieder E.  
373 2011. Predicting antigenic sites on the foot-and-mouth disease virus capsid of the South  
374 African Territories types using virus neutralization data. *J Gen Virol* 92:2297–2309.
- 375 19. Tesfaye Y, Khan F, Yami M, Wadsworth J, Knowles NJ, King DP, Gelaye E. 2020. A  
376 vaccine-matching assessment of different genetic variants of serotype O foot-and-mouth  
377 disease virus isolated in Ethiopia between 2011 and 2014. *Arch Virol* 165:1749–1757.
- 378 20. Upadhyaya S, Mahapatra M, Mioulet V, Parida S. 2021. Molecular Basis of Antigenic  
379 Drift in Serotype O Foot-and-Mouth Disease Viruses (2013-2018) from Southeast Asia.  
380 *Viruses* 13.
- 381 21. Yang M, Xu W, Goolia M, Zhang Z. 2014. Characterization of monoclonal antibodies  
382 against foot-and-mouth disease virus serotype O and application in identification of  
383 antigenic variation in relation to vaccine strain selection. *Virol J* 11:136.
- 384 22. World Organization for Animal Health. 2022. Foot and Mouth Disease (Infection with  
385 Foot and Mouth Disease Virus), p. 1–34. *In* OIE Terrestrial Manual.
- 386 23. Larsson A. 2014. AliView: a fast and lightweight alignment viewer and editor for large  
387 datasets. *Bioinformatics* 30:3276–3278.
- 388 24. Kumar S, Stecher G, Li M, Knyaz C, Tamura K. 2018. MEGA X: Molecular Evolutionary  
389 Genetics Analysis across Computing Platforms. *Mol Biol Evol* 35:1547–1549.
- 390 25. Zuckerkandl E, Pauling L. 1965. Evolutionary Divergence and Convergence in Proteins,  
391 p. 97–166. *In* *Evolving Genes and Proteins*. Elsevier.
- 392 26. CRAN - Package stringr. <https://stringr.tidyverse.org/authors.html>. Retrieved 17 June  
393 2022.
- 394 27. Keck F. 2020. Handling biological sequences in R with the bioseq package. *Methods Ecol*  
395 *Evol* 11:1728–1732.
- 396 28. Paradis E, Schliep K. 2019. ape 5.0: an environment for modern phylogenetics and

- 397 evolutionary analyses in R. *Bioinformatics* 35:526–528.
- 398 29. CRAN - Package tidysq. <https://cran.rstudio.com/web/packages/tidysq/index.html>.  
399 Retrieved 17 June 2022.
- 400 30. R Core Team. 2022. R: A Language and Environment for Statistical Computing. Vienna,  
401 Austria.
- 402 31. Chirkova Z V., Filimonov SI, Prituzhalov I V., Vasanov EA, Abramov IG. 2017.  
403 Synthesis of chalcones from 3-formyl-substituted pyrrolo[3,4-f]indole-5,7-diones. *Russ*  
404 *Chem Bull* 66:882–885.
- 405 32. Paploski IAD, Makau DN, Pamornchainavakul N, Baker JP, Schroeder D, Rovira A,  
406 VanderWaal K. 2022. Potential Novel N-Glycosylation Patterns Associated with the  
407 Emergence of New Genetic Variants of PRRSV-2 in the U.S. *Vaccines* 10:2021.
- 408 33. Yao Y, Li X, Liao B, Huang L, He P, Wang F, Yang J, Sun H, Zhao Y, Yang J. 2017.  
409 Predicting influenza antigenicity from Hemagglutinin sequence data based on a joint  
410 random forest method. *Sci Rep* 7:1545.
- 411 34. Bedford T, Suchard MA, Lemey P, Dudas G, Gregory V, Hay AJ, McCauley JW, Russell  
412 CA, Smith DJ, Rambaut A. 2014. Integrating influenza antigenic dynamics with  
413 molecular evolution. *Elife* 2014.
- 414 35. Kuhn M. 2008. Building Predictive Models in R Using the caret Package. *J Stat Softw*  
415 28:1–26.
- 416 36. Liaw A, Wiener M. 2002. Classification and Regression by randomForest. *R News* 2:18–  
417 22.
- 418 37. Hastie T, Tibshirani R, Friedman J. 2009. Random Forests, p. 587–604. *In* Springer (ed.),  
419 *The Elements of Statistical Learning* Second. Springer Series in Statistics, New York, NY.
- 420 38. Machado G, Mendoza MR, Corbellini LG. 2015. What variables are important in  
421 predicting bovine viral diarrhea virus? A random forest approach. *Vet Res* 46:1–15.
- 422 39. Makau DN, Prieto C, Martínez-Lobo FJ, Paploski IAD, VanderWaal K. 2023. Predicting  
423 Antigenic Distance from Genetic Data for PRRSV-Type 1: Applications of Machine  
424 Learning. *Microbiol Spectr* 11.
- 425 40. Bachanek-Bankowska K, Wadsworth J, Henry E, Ludi AB, Bin-Tarif A, Statham B, King  
426 DP, Afzal M, Hussain M, Manzoor S, Abubakar M, Knowles NJ. 2019. Genome  
427 Sequences of Antigenically Distinct Serotype O Foot-and-Mouth Disease Viruses from  
428 Pakistan. *Microbiol Resour Announc* 8.
- 429 41. Neher RA, Russell CA, Shraiman BI. 2014. Predicting evolution from the shape of  
430 genealogical trees. *Elife* 3.
- 431 42. Xia Y-L, Li W, Li Y, Ji X-L, Fu Y-X, Liu S-Q. 2021. A Deep Learning Approach for  
432 Predicting Antigenic Variation of Influenza A H3N2. *Comput Math Methods Med*  
433 2021:1–10.
- 434 43. Forghani M, Khachay M. 2020. Convolutional Neural Network Based Approach to In

- 435 Silico Non-Anticipating Prediction of Antigenic Distance for Influenza Virus. *Viruses*  
436 12:1019.
- 437 44. Lee EK, Tian H, Nakaya HI. 2020. Antigenicity prediction and vaccine recommendation  
438 of human influenza virus A (H3N2) using convolutional neural networks. *Hum Vaccin*  
439 *Immunother* 16:2690–2708.
- 440 45. Paton DJ, Reeve R, Capozzo A V., Ludi A. 2019. Estimating the protection afforded by  
441 foot-and-mouth disease vaccines in the laboratory. *Vaccine* 37:5515–5524.
- 442 46. Lavoria M ángeles, Di-Giacomo S, Bucafusco D, Franco-Mahecha OL, Pérez-Filgueira  
443 DM, Capozzo AV. 2012. Avidity and subtyping of specific antibodies applied to the  
444 indirect assessment of heterologous protection against Foot-and-Mouth Disease Virus in  
445 cattle. *Vaccine* 30:6845–6850.
- 446 47. Robiolo B, La Torre J, Maradei E, Beascochea CP, Perez A, Seki C, Smitsaart E,  
447 Fondevila N, Palma E, Goris N, De Clercq K, Mattion N. 2010. Confidence in indirect  
448 assessment of foot-and-mouth disease vaccine potency and vaccine matching carried out  
449 by liquid phase ELISA and virus neutralization tests. *Vaccine* 28:6235–6241.
- 450 48. Mattion N, Goris N, Willems T, Robiolo B, Maradei E, Beascochea CP, Perez A,  
451 Smitsaart E, Fondevila N, Palma E, De Clercq K, La Torre J. 2009. Some guidelines for  
452 determining foot-and-mouth disease vaccine strain matching by serology. *Vaccine*  
453 27:741–747.
- 454 49. Makau DN, Lycett S, Michalska-Smith M, Paploski IAD, Cheeran MC-J, Craft ME, Kao  
455 RR, Schroeder DC, Doeschl-Wilson A, VanderWaal K. 2022. Ecological and evolutionary  
456 dynamics of multi-strain RNA viruses. *Nat Ecol Evol* 1–9.
- 457 50. Islam MR, Rahman MS, Amin M Al, Alam ASMRU, Siddique MA, Sultana M, Hossain  
458 MA. 2021. Evidence of combined effect of amino acid substitutions within G-H and B-C  
459 loops of VP1 conferring serological heterogeneity in foot-and-mouth disease virus  
460 serotype A. *Transbound Emerg Dis* 68:375–384.
- 461
- 462