1

# Spatio-temporal Dynamics of Foot-and-Mouth Disease Virus in South America

Luiz Max Carvalho[1*], Nuno Rodrigues Faria[2], Guido König[3], Marc A. Suchard[4,5], Philippe Lemey[6], Waldemir de Castro Silveira[7], Guy Baele[6]

**1 School of Applied Mathematics, Getúlio Vargas Foundation, Rio de Janeiro, Brazil.**

**2 Department of Zoology, University of Oxford, Oxford, United Kingdom.**

**3 Institute of Agrobiotechnology and Molecular Biology, INTA-CONICET, Buenos Aires, Argentina.**

**4 Departments of Biomathematics and Human Genetics, David Geffen School of Medicine at UCLA, University of California, Los Angeles, United States of America.**

**5 Department of Biostatistics, UCLA Fielding School of Public Health, University of California, Los Angeles, United States of America.**

**6 Department of Microbiology, Immunology and Transplantation, Rega Institute, KU Leuven, Leuven, Belgium.**

**7 Research and Development Division, Trimatrix Applied Biotechnology Ltd, Rio de Janeiro, Brazil.**

∗ **E-mail: luiz.fagundes@fgv.br**

## Abstract

Despite the decrease in incidence of foot-and-mouth disease virus (FMDV) in South America over the last years, the pathogen still circulates in the region and the risk of re-emergence in previously FMDV-free areas is a veterinary public health concern. In this paper, we employ modern phylodynamic methods to merge epidemiological and genetic data and reconstruct spatiotemporal patterns and determinants of the dispersal of the two most prevalent FMDV serotypes A and O in South America, while accounting for temporal and spatial sampling bias. We find that serotypes A and O differ in their temporal pattern of population dynamics, with serotype A displaying more temporal oscillation. Spatially, we traced the origins of the 2011 Paraguay outbreak to Argentina (posterior probability 0.5) and Brazil (posterior probability 0.36). Overall, we find that FMDV spread seems to happen mostly through transnational borders, with few long range transmissions, such as a well-supported link between Argentina and Peru for serotype A. We also found the trade of different livestock (pigs for serotype A and cattle for serotype O) to be associated with viral spread, providing a possible explanation for this pattern of spread. Our results showcase the usefulness of phylodynamic methods to the study and surveillance of FMDV in South America.

Key-words: foot-and-mouth disease virus, South America, animal trade, pathogen phylodynamics, phylogenetics, Bayesian inference, BEAST.

## Introduction

Foot-and-mouth disease virus (FMDV) is a rapidly evolving picornavirus and the causative agent of foot-and-mouth disease (FMD), the most important disease of domestic and wild cloven-hoofed animals (Grubman and Baxt, 2004). The virus can be classified in seven serotypes, three of which (A, O, and C) have circulated in South America. Serotype A caused large epidemics throughout the Southern cone in recent years (Perez et al., 2008; Malirat et al., 2012), while endemic circulation has been mostly limited to Venezuela (Malirat et al., 2012). Historically, serotype O has been the most prevalent serotype on the continent, and was limited to areas in the Andean region, in particular to Ecuador (Malirat et al., 2011), until recent outbreaks in Colombia (OiE, 2017, 2018). Serotype C on the other hand was last encountered in the continent in 1995 in Brazil (Correa Melo et al., 2002). Historical reports suggest

2

that FMDV arrived in South America in the late years of the 19th century with European colonization (Naranjo and Cosivi, 2013; Tully and Fares, 2008). By the 1970s, FMD was widespread in the region, with several large-scale epidemics being caused by multiple subtypes (Saraiva, 2003). In South America, FMD control and eradication has traditionally been pursued using a combination of mass vaccination programs (Saraiva and Darsie, 2004) and control of animal movements from areas in which FMDV infection was suspected. Over time, passive and active surveillance programs have, with different degrees of success, managed the early detection of FMDV. In order to achieve complete eradication, however, the strains involved in epidemics - especially those in previously FMDV-free areas - need to be accurately characterised and their temporal and spatial dyamics, studied.

Phylogenetic analyses have proven useful in recovering the transmission pathways from genetic data (Cottam et al., 2008b,a) and in providing insight into the processes that drive re-emergence (Di Nardo et al., 2011). More recently, molecular epidemiology tools have been used to infer the origin and evolutionary history of emerging strains in South America (Perez et al., 2008; Malirat et al., 2007, 2011; Maradei et al., 2013). However, as pointed out by Di Nardo et al. (2011), a common feature of FMDV molecular epidemiology studies is that joint evaluation of epidemiological, environmental and genetic data has usually been performed outside of a unified quantitative framework. The link between phylogenetic analyses and population- and host-specific factors – such as animal trade and vaccination – is usually established in a post-hoc rather than a model-based fashion. In the face of many sources of information, ranging from genetic data to environmental data on host distribution and outbreak counts, it's desirable to have a framework capable of integrating these sources of information coherently (Lemey et al., 2014; Dudas et al., 2017).

The field of (pathogen) phylodynamics combines population genetics and epidemiology to explicitly model the interaction between ecological processes such as migration and selection and the shape of the phylogenies (Grenfell et al., 2004; Volz et al., 2013; Dudas et al., 2017). Bayesian phylodynamics offers an attractive statistical framework to combine multiple sources of information while marginalizing over the topology space, thus accommodating phylogenetic uncertainty. In particular, phylogeographic methods can be employed to understand viral spatial dynamics under explicit spatial diffusion models (Lemey et al., 2009). Further, an important research goal is to gain insight into the major determinants of FMDV spread in the continent. Since human and animal movements constitute a major threat to eradication programs (Schley et al., 2009), using animal trade data as predictors can be a valuable tool to understand the role of livestock commerce in the spread of FMDV. For example, Nelson et al. (2011) coupled swine trade data and genetic data to show that swine movements in the United States drove the spread of a novel influenza virus of the H1 subtype while Lemey et al. (2014) used air travel data to study the driving factors of Influenza H3N2 spread across the globe.

Here, we investigate the phylodynamic patterns of serotypes A and O in South America using all publicly available VP1 (1D) sequences for those serotypes in South America, sampled over a long time-period (1955-2010 for serotype A and 1994-2011 for serotype O) in nearly all south American countries affected by FMD. We apply Bayesian phylodynamic methods to investigate the evolutionary dynamics of serotypes A and O in South America incorporating genetic, spatial and epidemiological data such as livestock trade, geographic distances and vaccination coverage. This flexible Bayesian phylodynamic framework allows for the testing of hypotheses concerning viral dispersal, while naturally accommodating phylogenetic uncertainty (Lemey et al., 2009; Faria et al., 2011; Lemey et al., 2014; Gill et al., 2016). We use BEAST (Suchard et al., 2018) to infer time-structured phylogenies and reconstruct past population dynamics, to which we overlay vaccination and serotype-specific notification data. To study the factors driving re-emergence, we use data on livestock trade and geographical distances as predictors for viral spatial diffusion and compare competing spatial dynamics models involving each predictor using a generalised linear model approach (Faria et al., 2013; Lemey et al., 2014). In addtion, we reconstruct the past population dynamics of both serotypes using a sampling-aware method (Karcher et al., 2016, 2020) that takes into account possible temporal sampling bias in the collection of sequences.

# Results

We searched GenBank for FMDV sequences, filtering those that contained the 1D (VP1) gene (over 6,900 sequences) and then keeping those that had location and year of isolation and belonged to South America (see Methods), resulting in final data sets of 184 (1955 − 2013) sequences for serotype A and 210 (1958 − 2011) sequences for serotype O. The maximum clade credibility (MCC) phylogenetic trees shown in Figure 1 point to a considerable amount of geographic movement, as indicated by the interspersing of sequences of different countries. The tree for serotype A (Figure 1A) shows two major clades that diverge early on, one containing most of the Argentinian sequences, the other clade being more geographically heterogeneous. For serotype O (Figure 1B), we also notice an early split, with a long-lasting Venezuelan lineage that persisted into the 1970s. A major lineage containing all of the Ecuadorian sequences is interspersed with sequences from Colombia. The time of the most recent common ancestor (tMRCA) for serotype A was estimated at 1919 (1903-1935), while the tMRCA for serotype O was estimated at 1946 (1924–1966), which covers the date of early sequences from 1958 (excluded from the analyses, see methods). Relaxed molecular clock analysis showed both serotypes have high substitution rates with the rate estimated at 3.52 (2.80–4.27) $\times 10^{-3}$ substitutions/site/year (s/s/y) for serotype A and at 5.41 (4.00–6.92) $\times 10^{-3}$ s/s/y for serotype O. Both serotypes showed considerable inter-lineage rate variation with coefficients of variation (posterior mean) of 1 and 1.6 for serotypes A and O, respectively.

Reconstructions of past population dynamics for both serotypes under naive and sampling-aware models are presented in Figure 2. The sampling-aware models account for dependence of the sampling process on the underlying effective population size ($N_e(t)$) and other simple time-varying covariates, allowing one to assess the presence of bias. Results suggest that effective population sizes for serotype A display a pattern of steady increase until circa the 1970s and then steady decline, which then becomes faster closer to the present. For serotype O, the naive reconstruction, which does not take preferential sampling into account, shows substantial oscillations, not present in the reconstructions using the preferential sampling model. Reconstructions that do not account for preferential sampling also lead to considerably wider Bayesian credibility intervals. We used MCC trees obtained from three independent chains per serotype as fixed phylogenies for the population dynamics reconstructions, and results were largely consistent across these replicates. In particular, model 3 ($\{\gamma(t), -t, -t^2\}$, see Text S1), which includes both a linear and a quadratic term on the sampling time, yielded the highest (log) marginal likelihood for all three replicates of serotype A and for two of three replicates for serotype O. For all of the concordant replicates, the smallest log-Bayes factors in support of model 3 were ≈ 9 for serotype A and ≈ 3 for serotype O. This less decisive support in favour of model 3 for serotype O is also manifested in the discrepant replicate, in which model 2, which includes only a linear term on $t$, is favoured with a log Bayes factor of ≈ 1.7.

Figure 3 shows the network of FMDV spread for both serotypes in South America, reconstructed using Bayesian stochastic search variable selection (BSSVS, Lemey et al. (2009)) and stochastic mapping (Minin and Suchard, 2008). Many connections, for example the Argentina-Brazil, Argentina-Uruguay, Venezuela-Colombia and Brazil-Venezuela links are shared between serotypes, with varying degrees of statistical support. Overall, connections exist mostly between countries that share borders, and these cross-border connections are mostly concordant across serotypes. A notable exception is a strong connection between Argentina and Peru for serotype A, which is absent in the network for serotype O. Peru seems to be a hub for serotype A, with most connections being imports into the country rather than radiating out of it. When we look at source-sink dynamics by computing the net exchange rate of a country, i.e., the difference between the expected transitions from and to that country, differences between serotypes emerge. For instance, Brazil is a sink for serotype A, but acts as a source for serotype O. The opposite is true for Colombia. On the other hand, Ecuador is a source and a sink (i.e. a hub) for serotype A, while Argentina seems to act as both source and a sink for serotype O. In addition, while the largest positive net exchange rates (imports minus exports) are similar for both serotypes, negative exchange rates differ. Argentina and Colombia are strong sinks for serotype O, but we do not observe a similar pattern for serotype A.

4

We traced the origins of specific epidemics and found that the origin of the 2001 serotype A lies within the root of the tree (Figure 4A), which indicates that the strains responsible for the 2001 epidemic were likely already present in the country – which is the most likely state of the root for serotype A. We traced the MRCA of the 2002 serotype A Ecuadorian epidemic to Brazil with posterior probability 0.52, while Peru (0.2) and Venezuela (0.17) were alternative locations of origin (Figure 4B). For serotype O, we also traced the origins of the 2002 epidemic in Ecuador, and found overwhelming evidence of a Colombian origin (posterior probability 0.92, Figure 4C). We estimated the most probable location of origin of the serotype O Paraguay 2011 isolate (Figure 4D) and found that Argentina is the most likely place of origin (posterior probability 0.5) followed by Brazil (0.36, see Discussion). Tracing the origins of the serotype A 2001 epidemics in Brazil and Uruguay revealed multiple countries of origin (Figure S8A and B), while the Bolivian 2001 strains probably originated in Peru (posterior probability (0.47, Figure S8C). The 2008 strains of serotype A in Colombia and serotype O in Peru can be traced with almost complete certainty to Venezuela and Ecuador, respectively (Figure S8D and F) – see Lineage 1 in Malirat et al. (2012).

To evaluate possible drivers of FDMV spread in South America, we employed the generalised linear model framework of Lemey et al. (2014). The results of this spatial GLM analysis are summarised in Figure 5, and show that out of 15 predictors, only the trade of pigs (1995 − 2004) and cattle (1986 − 1994) are significant predictors of spread for serotypes A and O respectively. However, other predictors that failed to be included with substantial probability, nevertheless yielded posterior 95% BCIs for the coefficients that exclude zero. Examples include the product of the number of sequences and pigs trade (2004−2013) for serotype A, and the presence of borders, pigs trade (1995−2004), product of the number of sequences and number of sequences from the destination location for serotype O (Figure 5B).

## Discussion

Dating analyses showed serotype A to be older than serotype O in the continent, a finding consistent with a global analysis of both serotypes (Tully and Fares, 2008). Serotype O also showed remarkably higher substitution rates, but caution should be exercised when interpreting these findings, as the substitution rate is a quantity that incorporates molecular characteristics linked to replication as well as populational processes driving mutation fixation (Holmes et al., 2016). Our population dynamics reconstructions suggest that serotype A has experienced a rise and fall trajectory with a sharp increase and subsequent decrease of its effective population size. Serotype O, in contrast, shows a long period of stability in terms of its population size, followed by a period of decline starting in the mid 1990s. Overall, effective population sizes follow the downward trend of FMD case counts shown in Figure S7, but the limited temporal span of the case data (1990 to 2010) precludes more general conclusions. Naranjo and Cosivi (2013) provide a longer time series (Figure 3 therein) which shows a spike in cases for serotype A around the mid 1970s and that serotype O cases show a decreasing trend from the 1970s which intensifies around the mid 1990s, which coincides with the start of the decline seen in the left panel of Figure 2 in the present paper. Interestingly, the 2001 serotype A epidemic does not seem to lead to noticeable changes in diversity.

FMDV spread in South America seems to happen mostly through transnational borders, with long migration routes being rarer. In keeping with the spatial Bayes factors results (Figure 3), we found that the presence of borders between countries attained a positive albeit not significant (BF < 3) coefficient (Figure 5). This is specially true for serotype O, for which the 95% BCI for the coefficient excluded zero (Figure 5B). These results are consistent also with the fact that geographic distance yielded a negative coefficient, with the 95% BCI for serotype O excluding zero. Taken together, these results indicate that, especially for serotype O, FMDV dispersal takes place over shorter distances, across shared borders and is connected with the trade of live pigs and cattle. The finding that the Ecuadorian epidemic was most likely seeded from Colombia (Figure 4C) also suggests an effect of cattle trade, since Carvalho et al. (2013) show that the origin of the Ecuadorian epidemic was the province of Esmeraldas at the border between

the two countries and this is where most of the Ecuadorian cattle trade takes place (Maradei et al., 2011). Distance to the border is considered one of the most important factors to FMD spread, at least in the context of the Paraguay/Brazil border (Amaral et al., 2016). On the other hand, our finding that Argentina was the most likely origin of the 2011 Paraguayan isolate is in disagreement with Maradei et al. (2013), who state that, antigenically, the O/San Pedro/Par/11 isolate is not very similar to Argentinian strains such as O/Corrientes/Arg/06. We note that our analysis here takes genomic information into account, not just phenotypic traits, which might explain the disagreement.

Figure S8F shows that Ecuador is most certainly the origin of the Peruvian 2008 outbreak, which may be connected with the illegal movement of cattle across the border, as suggested by Correa Melo et al. (2002). Regarding the reservoirs of FMDV in the continent, our results for serotype O partially confirm common knowledge (e.g.(Saraiva, 2003) and (Naranjo and Cosivi, 2013)) that the (possibly illegal) trade of cattle is a driving force of FMDV spread. On the other hand, our results for serotype A seem to implicate the swine population in the spread of the virus, a mechanism that could be facilitated by swine and caprines only being vaccinated during public health emergencies (Saraiva, 2003).

A crucial goal of the present study was to perform phylodynamic inference while accounting for preferential sampling, both in time and space. In order to tackle temporal preferential sampling, we have employed a model that explicitly incorporates the sequence sampling dates (Karcher et al., 2020), while also allowing for the inclusion of covariates to help increase explanatory power. Our finding that accounting for preferential sampling leads to narrower credibility intervals and better fitting models highlights the dangers of "naive" reconstructions based on non-uniform temporal sampling. In particular, we found that our sampling-aware reconstructions showed fewer oscillations, which could be mistakenly interpreted as genuine changes in population dynamics under a naive model. In Text S1, we provide further simulations that investigate the effect of the specific temporal sampling patterns observed for both serotypes when compared to uniform temporal sampling.

To account for spatial preferential sampling, we included the number of sequences isolated in each location as a predictor of spread. Since our goal was a conservative analysis, we included the sampling structure in different ways – difference and product of the number of sequences between countries, numbers of sequences in both origin and destination – in order to maximise our ability to detect an effect of sampling. Our results show a mild, not statistically significant effect of sampling as evidenced by the zero-excluding 95% BCIs obtained for the predictors "Product # sequences" and " # sequences destination" for serotype O. In addition to *detecting* spatial preferential sampling, the inclusion of sampling-related predictors allows for *accommodating* preferential sampling so that other coefficient estimates, made conditionally, are robust to it.

Finally, it is important to point out that differences between serotypes should not be readily attributed to inherent biological differences. Most inferences reported here depend on complex population-level processes such as the migration of infected hosts, immunological composition of host populations and rearing practices. As an example, we note that even quantities more strongly linked to molecular processes such as the evolutionary rate depend on population-level processes that obscure or magnify biological differences in mutation rate by influencing rates of fixation.

## Limitations of this study

Sampling-related effects have the potential to lead to incorrect inferences about temporal and spatial patterns in phylodynamic studies (Hall et al., 2016; Dearlove et al., 2017). As discussed above, we have employed modern statistical methods to quantify and mitigate spatial and temporal preferential sampling in the data analysed in this study. Despite this, however, we feel it is important to clarify its limitations as they relate to the nature of the sampling process. We highlight three key limitations, which in our view cannot be completely circumvented by the use of state-of-the-art statistical methods. First, it is important to understand that while the methods employed in this paper help quantify and account for the effects of preferential sampling, ultimately biased data lack information and hence force us to draw

6

more conservative conclusions. This is particularly important to consider in light of the fact that early reporting for FMD was not reliable and sampling is likely to be significantly biased downward in the period before 1980s. Secondly, the sequences in the present study are from a particular gene (1D, VP1) rather than the full genomes, for which there are considerably fewer sequences that originate only from a few countries. This restricts our ability to resolve finer details in the phylogenies and by extension limits the ability of phylodynamic models (Valdazo-González et al., 2012; Dudas and Bedford, 2019). A final difficulty we would like to point out is the lack of detailed covariate information, especially negative controls (see Dellicour et al. (2018) for a description of negative controls in the context of phylogenetic GLMs). For example, due to the fact that our phylogenies spanned periods of time much longer than the available data, we could not include any temporal predictors of FMDV dynamics. In the supplementary materials, we provide all of the data collected for this study, including data that we could not employ directly, such as vaccination and animal production data. We hope these data can be of use to other researchers in future studies of more limited temporal and spatial scope.

## Conclusions

In this paper, we have assembled a large data set of FMDV VP1 gene sequences from South America and employed state-of-the-art statistical methods to analyse the spatiotemporal dynamics of the virus, while accounting for preferential sampling. We found that FMDV spread occurs mostly through shared borders and is connected to livestock trade. We have traced the origins of several epidemics on the continent, providing valuable information to policy makers. Ultimately, however, the convenience nature of the sampling process creates biases that cannot be easily accommodated, a problem which is made worse by the sequencing of genes rather than full genomes (Dudas and Bedford, 2019). Nevertheless, our study showcases the power of phylodynamic methods in aiding epidemiological surveillance of FMD in South America and prompts public health agencies to adopt sequencing as an integral part of their surveillance apparatus.

## Methods

In this section we detail the computational techniques employed to collect and analyse the data in this paper. We take a Bayesian approach to parameter inference and statistical testing and all parameter estimates reported in this paper are of the form: posterior mean (95% credibility interval), unless otherwise stated.

### Genetic and epidemiological data

#### Genetic data.

We retrieved all FMDV nucleotide sequences available from GenBank (Benson et al., 2013) from the National Center for Biotechnology Information (NCBI, http://www.ncbi.nlm.nih.gov/) with more than 600 bp. This first step yielded 6,907 sequences which were then filtered to exclude all sequences that did not include the 1D (VP1) gene, resulting in 4,507 sequences being kept. We then filtered for sequences from serotypes A and O, yielding 1051 and 2350 sequences, respectively. Next, we excluded sequences that had been extensively passaged in cell culture and selected all sequences from South America (Argentina, Bolivia, Brazil, Colombia, Ecuador, Paraguay, Peru, Uruguay, Venezuela) for which information on country and year of isolation was available.

This resulted in 185 sequences (from eight countries) for serotype A and 215 sequences (from nine countries) for serotype O, covering time spans of 62 (1951-2013) and 53 (1958-2011) years, respectively (see Tables S3 and S4 for details). We aligned each data set using the MAFFT (Katoh et al., 2002) algorithm implemented in the Geneious (Kearse et al., 2012) software package. After a preliminary phylogenetic

7

analysis (see below), an early serotype A sequence (Venezuela 1951) was excluded because its root-to-tip divergence was incompatible with its sampling date (Rambaut et al., 2016). For serotype O, five sequences were excluded under the same criteria (Table S1). The final data comprised 184 $(1955 - 2013)$ and 210 $(1958 - 2011)$ sequences for serotypes A and O, respectively.

**Acquisition of trade data.**

Data on animal trade were obtained from the FAO database (`http://faostat.fao.org/`). We retrieved data on the *detailed trade matrix* for cattle, pigs and sheep (number of live animals exchanged) covering the period from 1986 to 2009, for each of the nine countries. The detailed trade matrix includes information on export quantity, export value, import quantity and import value for each variable. We used the information on export and import quantity to compose our predictors. We refer to Supplementary Text S1 for more details.

**Vaccination and case data.**

Data on the number of vaccine doses (irrespective of serotype) from 1990 to 2010 were obtained from the PANAFTOSA annual reports (`https://www.paho.org/panaftosa/`) and then standardised to doses per (cattle) head in each country. Serotype-specific outbreak notifications were obtained from the FMD Bioportal (`http://fmdbioportal.ucdavis.edu:8080/`).

**Data availability.** All the data used in this paper including BEAST XML files are hosted at `https://github.com/maxbiostat/FMDV_AMERICA`.

In this study we take a Bayesian approach to testing evolutionary hypotheses while accommodating phylogenetic uncertainty. Details on computational settings and prior choices can be found in Text S1.

## Phylogenetic Analysis

We assume a general time reversible (GTR) (Tavaré, 1986) model of sequence evolution, along with gamma-distributed rate heterogeneity (4 categories) for all our analyses. We conducted an initial analysis of the two data sets described, employing PhyML (Guindon and Gascuel, 2003) to obtain maximum likelihood phylogenies which we then used in conjunction with TempEst (Rambaut et al., 2016) to produce root-to-tip divergence (RDV) plots and identify discrepant sequences (see Rambaut et al. (2016) for details). Having identified considerable rate variation among branches we used the Bayesian Evolutionary Analysis by Sampling Trees (BEAST) (Suchard et al., 2018) software package to infer time-structured phylogenies under relaxed clock models (Drummond et al., 2006), employing the BEAGLE (Ayres et al., 2012, 2019) library to gain computational efficiency. In order to select an appropriate combination of sequence evolution (Markov), coalescent and relaxed clock models for each data set (see Text S1), we use state-of-the-art marginal likelihood estimators, namely generalised stepping-stone sampling (SS) (Baele et al., 2016) implemented in BEAST (Suchard et al., 2018). Following a preliminary skygrid analysis (Gill et al., 2013), we selected the constant population size coalescent parametric model – and its associated working prior – as the tree prior for these comparisons.

## Spatio-temporal Dynamics

**Effective population size reconstruction**.

For an initial analysis, we employed the skygrid coalescent model (Gill et al., 2013; Hill and Baele, 2019) to reconstruct the past population dynamics of both serotypes. The skygrid model necessitates the specification of a cutoff value, $K$, commensurate with the age of the root (time to the most recent common ancestor, tMRCA) . We used $K = 150$ years for serotype A and $K = 100$ for serotype O sets as conservative estimates of the tMRCA of the sampled sequences. This demographic prior was used in all subsequent BEAST analysis.

In order to accommodate the effect of preferential (temporal) sampling, we employed the Bayesian Nonparametric Population Reconstruction under Preferential Sampling (BNPR-PS) model of Karcher et al. (2020) (see also Karcher et al. (2016)), that accounts for the sampling pattern of sequences through

8

an inhomogeneous Poisson process and allows the inclusion of time-varying covariates. We employed the function `BNPR_PS()` in the **phylodyn** (Karcher et al., 2017) R package to fit the preferential-sampling model using three MCC trees derived from previous BEAST analyses. Both the BNPR and skygrid models necessitate the estimation of a precision parameter, to which one usually assigns a Gamma distribution with parameters $\alpha = \beta = 0.001$. In this study however, we employed the penalised complexity (PC) prior of Simpson et al. (2017), which has better theoretical properties. More details on model formulation, testing and comparison are given in Text S1.

**Reconstructing geographical movements using BSSVS and stochastic mapping**

We employ Bayesian stochastic search variable selection (BSSVS) (Lemey et al., 2009) to compute Bayes factors for migration routes between countries and establish which routes are relevant to the spatial spread of FMDV in South America. Following the approach of Hall et al. (2013), we employ the stochastic mapping technique (Markov jump) of Minin and Suchard (2008) to obtain a probability distribution for the location (country) of origin of each of epidemic. To represent each epidemic (e.g. Ecuador 2002), we selected a group of representative sequences from that epidemic (country and year) and traced their origins.(see Text S1 for details)

**Phylogeographic generalised linear models**

In order to study the influence of different epidemiological predictors on viral diffusion through space, we used information on the trade of live cattle, pigs and sheep divided in three periods $(1986-1995, 1996-2004, 2005-2013)$ for each livestock, resulting in nine predictors. We also included the great circle distance between the (centroids of) countries and the presence/absence of borders between countries as predictors. Additionally, to assess the influence of sequence (spatial) sampling, we also included the product and absolute differences in sequence (samples) numbers between locations as a predictor of flow (see the appendix in Lemey et al. (2014) for a discussion). The numbers of sequences in each location were also included as origin and destination predictors (see Dudas et al. (2017)), leading to a total of 14 predictors being considered. In order to test the relevance of each predictor to spatial spread, we used Bayesian stochastic variable selection (BSSVS, Lemey et al. (2009)). All predictors with the exception of the borders were standard log-transformed. The relevance of each predictor can then be determined using Bayes factors (Kass and Raftery, 1995; Lemey et al., 2009, 2014). Details are given in Text S1.

## Software versions and computer programs

- **PhyML** version 3.0, downloaded from `www.atgc-montpellier.fr/phyml/`;

- **Tempest** version 1.5.1, downloaded from `http://tree.bio.ed.ac.uk/software/tempest/` ;

- **BEAST** version 1.10.4, downloaded from `https://github.com/beast-dev/beast-mcmc/releases/download/v1.10.4/`;

- **BEAGLE** version 3.0.0 downloaded from `https://github.com/beagle-dev/beagle-lib/releases`;

- **Jupyter** version 4.4.0, installed from the official Ubuntu repositories;

- **R** version 3.5.0, downloaded from `https://cran.r-project.org/`.

All analyses conducted within a GNU/Linux computational environment. Code to produce many of the plots/analyses is available from `https://github.com/maxbiostat/FMDV_AMERICA`.

## Acknowledgments

9

*Conflict of Interest:* none declared

10

# References

Amaral, T. B., Gond, V., and Tran, A. (2016). Mapping the likelihood of foot-and-mouth disease introduction along the border between brazil and paraguay. *Pesquisa Agropecuária Brasileira*, 51(5):661–670.

Ayres, D. L., Cummings, M. P., Baele, G., Darling, A. E., Lewis, P. O., Swofford, D. L., Huelsenbeck, J. P., Lemey, P., Rambaut, A., and Suchard, M. A. (2019). Beagle 3: Improved performance, scaling, and usability for a high-performance computing library for statistical phylogenetics. *Systematic biology*.

Ayres, D. L., Darling, A., Zwickl, D. J., Beerli, P., Holder, M. T., Lewis, P. O., Huelsenbeck, J. P., Ronquist, F., Swofford, D. L., Cummings, M. P., Rambaut, A., and Suchard, M. A. (2012). BEAGLE: an application programming interface and high-performance computing library for statistical phylogenetics. *Syst. Biol.*, 61(1):170–173.

Baele, G., Lemey, P., and Suchard, M. A. (2016). Genealogical working distributions for bayesian model testing with phylogenetic uncertainty. *Systematic biology*, 65(2):250–264.

Benson, D. A., Cavanaugh, M., Clark, K., Karsch-Mizrachi, I., Lipman, D. J., Ostell, J., and Sayers, E. W. (2013). GenBank. *Nucleic Acids Res.*, 41(Database issue):36–42.

Carvalho, L. M., Santos, L. B., Faria, N. R., and de Castro Silveira, W. (2013). Phylogeography of foot-and-mouth disease virus serotype O in Ecuador. *Infect. Genet. Evol.*, 13:76–88.

Correa Melo, E., Saraiva, V., and Astudillo, V. (2002). Review of the status of foot and mouth disease in countries of South America and approaches to control and eradication. *Rev. - Off. Int. Epizoot.*, 21(3):429–436.

Cottam, E. M., Thebaud, G., Wadsworth, J., Gloster, J., Mansley, L., Paton, D. J., King, D. P., and Haydon, D. T. (2008a). Integrating genetic and epidemiological data to determine transmission pathways of foot-and-mouth disease virus. *Proc. Biol. Sci.*, 275(1637):887–895.

Cottam, E. M., Wadsworth, J., Shaw, A. E., Rowlands, R. J., Goatley, L., Maan, S., Maan, N. S., Mertens, P. P., Ebert, K., Li, Y., Ryan, E. D., Juleff, N., Ferris, N. P., Wilesmith, J. W., Haydon, D. T., King, D. P., Paton, D. J., and Knowles, N. J. (2008b). Transmission pathways of foot-and-mouth disease virus in the United Kingdom in 2007. *PLoS Pathog.*, 4(4):e1000050.

Dearlove, B. L., Xiang, F., and Frost, S. D. (2017). Biased phylodynamic inferences from analysing clusters of viral sequences. *Virus evolution*, 3(2).

Dellicour, S., Vrancken, B., Trovão, N. S., Fargette, D., and Lemey, P. (2018). On the importance of negative controls in viral landscape phylogeography. *Virus evolution*, 4(2):vey023.

Di Nardo, A., Knowles, N. J., and Paton, D. J. (2011). Combining livestock trade patterns with phylogenetics to help understand the spread of foot and mouth disease in sub-Saharan Africa, the Middle East and Southeast Asia. *Rev. - Off. Int. Epizoot.*, 30(1):63–85.

Drummond, A. J., Ho, S. Y., Phillips, M. J., and Rambaut, A. (2006). Relaxed phylogenetics and dating with confidence. *PLoS biology*, 4(5):e88.

Dudas, G. and Bedford, T. (2019). The ability of single genes vs full genomes to resolve time and space in outbreak analysis. *bioRxiv*, page 582957.

Dudas, G., Carvalho, L. M., Bedford, T., Tatem, A. J., Baele, G., Faria, N. R., Park, D. J., Ladner, J. T., Arias, A., Asogun, D., et al. (2017). Virus genomes reveal factors that spread and sustained the ebola epidemic. *Nature*, 544(7650):309–315.

Faria, N. R., Suchard, M. A., Rambaut, A., and Lemey, P. (2011). Toward a quantitative understanding of viral phylogeography. *Curr Opin Virol*, 1(5):423–429.

Faria, N. R., Suchard, M. A., Rambaut, A., Streicker, D. G., and Lemey, P. (2013). Simultaneously reconstructing viral cross-species transmission history and identifying the underlying constraints. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 368(1614).

Gill, M. S., Lemey, P., Bennett, S. N., Biek, R., and Suchard, M. A. (2016). Understanding past population dynamics: Bayesian coalescent-based modeling with covariates. *Systematic biology*, 65(6):1041–1056.

Gill, M. S., Lemey, P., Faria, N. R., Rambaut, A., Shapiro, B., and Suchard, M. A. (2013). Improving bayesian population dynamics inference: a coalescent-based model for multiple loci. *Molecular biology and evolution*, 30(3):713–724.

Grenfell, B. T., Pybus, O. G., Gog, J. R., Wood, J. L., Daly, J. M., Mumford, J. A., and Holmes, E. C. (2004). Unifying the epidemiological and evolutionary dynamics of pathogens. *Science*, 303(5656):327–332.

Grubman, M. J. and Baxt, B. (2004). Foot-and-mouth disease. *Clin. Microbiol. Rev.*, 17(2):465–493.

Guindon, S. and Gascuel, O. (2003). A simple, fast, and accurate algorithm to estimate large phylogenies by maximum likelihood. *Syst. Biol.*, 52(5):696–704.

Hall, M. D., Knowles, N. J., Wadsworth, J., Rambaut, A., and Woolhouse, M. E. (2013). Reconstructing geographical movements and host species transitions of foot-and-mouth disease virus serotype sat 2. *MBio*, 4(5):e00591–13.

Hall, M. D., Woolhouse, M. E., and Rambaut, A. (2016). The effects of sampling strategy on the quality of reconstruction of viral population dynamics using bayesian skyline family coalescent methods: A simulation study. *Virus evolution*, 2(1).

Hill, V. and Baele, G. (2019). Bayesian estimation of past population dynamics in beast 1.10 using the skygrid coalescent model. *Molecular biology and evolution*.

Holmes, E. C., Dudas, G., Rambaut, A., and Andersen, K. G. (2016). The evolution of Ebola virus: Insights from the 2013–2016 epidemic. *Nature*, 538(7624):193–200.

Karcher, M., Palacios, J., Lan, S., and Minin, V. (2017). phylodyn: an R package for phylodynamic simulation and inference. *Molecular Ecology Resources*, 17(1):96–100.

Karcher, M. D., Carvalho, L. M., Suchard, M. A., Dudas, G., and Minin, V. N. (2020). Estimating effective population size changes from preferentially sampled genetic sequences. *PLoS Computational Biology*.

Karcher, M. D., Palacios, J. A., Bedford, T., Suchard, M. A., and Minin, V. N. (2016). Quantifying and mitigating the effect of preferential sampling on phylodynamic inference. *PLoS computational biology*, 12(3):e1004789.

Kass, R. and Raftery, A. (1995). Bayes Factors. *J. Amer. Statist. Assoc.*, 90(430):773–795.

Katoh, K., Misawa, K., Kuma, K.-i., and Miyata, T. (2002). Mafft: a novel method for rapid multiple sequence alignment based on fast fourier transform. *Nucleic acids research*, 30(14):3059–3066.

12

Kearse, M., Moir, R., Wilson, A., Stones-Havas, S., Cheung, M., Sturrock, S., Buxton, S., Cooper, A., Markowitz, S., Duran, C., et al. (2012). Geneious basic: an integrated and extendable desktop software platform for the organization and analysis of sequence data. *Bioinformatics*, 28(12):1647–1649.

Lemey, P., Rambaut, A., Bedford, T., Faria, N., Bielejec, F., Baele, G., Russell, C. A., Smith, D. J., Pybus, O. G., Brockmann, D., and Suchard, M. A. (2014). Unifying viral genetics and human transportation data to predict the global transmission dynamics of human influenza H3N2. *PLoS Pathog.*, 10(2):e1003932.

Lemey, P., Rambaut, A., Drummond, A. J., and Suchard, M. A. (2009). Bayesian phylogeography finds its roots. *PLoS Comput. Biol.*, 5:e1000520.

Malirat, V., Bergmann, I. E., Campos, R. d. e. M., Salgado, G., Sanchez, C., Conde, F., Quiroga, J. L., and Ortiz, S. (2011). Phylogenetic analysis of Foot-and-Mouth Disease Virus type O circulating in the Andean region of South America during 2002-2008. *Vet. Microbiol.*, 152:74–87.

Malirat, V., Bergmann, I. E., de Mendonca Campos, R., Conde, F., Quiroga, J. L., Villamil, M., Salgado, G., and Ortiz, S. (2012). Molecular epidemiology of foot-and-mouth disease virus type A in South America. *Vet. Microbiol.*, 158(1-2):82–94.

Malirat, V., de Barros, J. J., Bergmann, I. E., Campos, R. d. e. M., Neitzert, E., da Costa, E. V., da Silva, E. E., Falczuk, A. J., Pinheiro, D. S., de Vergara, N., Cirvera, J. L., Maradei, E., and Di Landro, R. (2007). Phylogenetic analysis of foot-and-mouth disease virus type O re-emerging in free areas of South America. *Virus Res.*, 124(1-2):22–28.

Maradei, E., Malirat, V., Beascoechea, C. P., Benitez, E. O., Pedemonte, A., Seki, C., Novo, S. G., Balette, C. I., D'Aloia, R., La Torre, J. L., Mattion, N., Toledo, J. R., and Bergmann, I. E. (2013). Characterization of a type O foot-and-mouth disease virus re-emerging in the year 2011 in free areas of the Southern Cone of South America and cross-protection studies with the vaccine strain in use in the region. *Vet. Microbiol.*, 162(2-4):479–490.

Maradei, E., Perez Beascoechea, C., Malirat, V., Salgado, G., Seki, C., Pedemonte, A., Bonastre, P., D'Aloia, R., La Torre, J. L., Mattion, N., Rodriguez Toledo, J., and Bergmann, I. E. (2011). Characterization of foot-and-mouth disease virus from outbreaks in Ecuador during 2009-2010 and cross-protection studies with the vaccine strain in use in the region. *Vaccine*, 29:8230–8240.

Minin, V. N. and Suchard, M. A. (2008). Counting labeled transitions in continuous-time Markov models of evolution. *J Math Biol*, 56(3):391–412.

Naranjo, J. and Cosivi, O. (2013). Elimination of foot-and-mouth disease in South America: lessons and challenges. *Philos. Trans. R. Soc. Lond., B, Biol. Sci.*, 368(1623):20120381.

Nelson, M. I., Lemey, P., Tan, Y., Vincent, A., Lam, T. T., Detmer, S., Viboud, C., Suchard, M. A., Rambaut, A., Holmes, E. C., and Gramer, M. (2011). Spatial dynamics of human-origin H1 influenza A virus in North American swine. *PLoS Pathog.*, 7(6):e1002077.

OiE (2017). Foot-and-mouth disease, colombia.

OiE (2018). Foot-and-mouth disease, colombia.

Perez, A. M., Konig, G., Spath, E., and Thurmond, M. C. (2008). Variation in the VP1 gene of foot-and-mouth disease virus serotype A associated with epidemiological characteristics of outbreaks in the 2001 epizootic in Argentina. *J. Vet. Diagn. Invest.*, 20(4):433–439.

Rambaut, A., Lam, T. T., Carvalho, L. M., and Pybus, O. G. (2016). Exploring the temporal structure of heterochronous sequences using tempest (formerly path-o-gen). *Virus Evolution*, 2(1):vew007.

Saraiva, V. (2003). Epidemiology of Foot-and-mouth disease in South America. In Dodet, B. and Vicari, M., editors, *Foot and mouth disease: control strategies*, pages 43–54. Paris: Elsevier SAS.

Saraiva, V. and Darsie, G. (2004). The use of vaccines in South American foot-and-mouth disease eradication programmes. *Dev Biol (Basel)*, 119:33–40.

Schley, D., Gubbins, S., and Paton, D. J. (2009). Quantifying the risk of localised animal movement bans for foot-and-mouth disease. *PLoS ONE*, 4(5):e5481.

Simpson, D., Rue, H., Riebler, A., Martins, T. G., Sørbye, S. H., et al. (2017). Penalising model component complexity: A principled, practical approach to constructing priors. *Statistical Science*, 32(1):1–28.

Suchard, M. A., Lemey, P., Baele, G., Ayres, D. L., Drummond, A. J., and Rambaut, A. (2018). Bayesian phylogenetic and phylodynamic data integration using beast 1.10. *Virus Evolution*, 4(1):vey016.

Tavaré, S. (1986). *Some Probabilistic and Statistical Problems in the Analysis of DNA Sequences*, volume 17, pages 57–86. Amer Mathematical Society.

Tully, D. C. and Fares, M. A. (2008). The tale of a modern animal plague: tracing the evolutionary history and determining the time-scale for foot and mouth disease virus. *Virology*, 382(2):250–256.

Valdazo-González, B., Polihronova, L., Alexandrov, T., Normann, P., Knowles, N. J., Hammond, J. M., Georgiev, G. K., Özyörük, F., Sumption, K. J., Belsham, G. J., et al. (2012). Reconstruction of the transmission history of rna virus outbreaks using full genome sequences: foot-and-mouth disease virus in bulgaria in 2011. *PLoS one*, 7(11):e49650.

Volz, E. M., Koelle, K., and Bedford, T. (2013). Viral phylodynamics. *PLoS Comput. Biol.*, 9(3):e1002947.

(A) Serotype A



(B) Serotype O

**Figure 1. Maximum clade credibility (MCC) phylogenetic trees reconstructed for FMDV serotypes A and O, with country information mapped**. Internal branches are coloured according to the most probable country. Please notice the slight difference in colour scheme in the legend of both panels.

**Figure 2. Reconstructed population dynamics for serotypes A and O in South America.** In the left panel we show reconstructions that do not take sampling (sequencing) patterns into account, whereas the right panel shows reconstructions using a preferential sampling model (Karcher et al., 2020) that includes a quadratic term on the log-intensity, $\{\gamma(t), -t, -t^2\}$ and which yielded the highest (log) marginal likelihood among the coalescent models tested (see Text S1 for more details). Colours depict serotypes.

16

A

O



Bayes factors
— under 12.14
— 12.14 – 25.19
▬ over 25.19

Net migration rate
☐ under −0.04
☐ −0.04 – 0.34
☐ 0.34 – 1.7
☐ over 1.7

Bayes factors
— under 11.37
— 11.37 – 21.1
▬ over 21.1

Net migration rate
☐ under −1.91
☐ −1.91 – −0.12
☐ −0.12 – 1.63
☐ over 1.63

**Figure 3. Migration networks for FMDV serotypes A and O in South America**. We estimate the number of migration events between countries using the stochastic mapping (Markov jumps) technique of Minin and Suchard (2008). Additionally, we perform a separate analysis employing Bayesian Stochastic Search Variable Selection (BSSVS) to determine the most significant migration routes. Bayes factors are depicted by arrows, with line thickness proportional to BF magnitude (we only plot BFs bigger than 3). Coroplethic maps show the net migration rates (from - to) for each country but colour scales vary by panel (serotype).

**A Argentina early 2001**

Origin of virus

Root — 0.56

Uruguay — 0.4

Other — 0.04

0.00  0.25  0.50  0.75  1.00
Posterior probability

(A)

**A Ecuador 2002**

Origin of virus

Brazil — 0.52

Peru — 0.2

Venezuela — 0.17

Argentina — 0.09

Other — 0.01

0.00  0.25  0.50  0.75  1.00
Posterior probability

(B)

**O Ecuador 2002**

Origin of virus

Colombia — 0.92

Other — 0.08

0.00  0.25  0.50  0.75  1.00
Posterior probability

(C)

**O Paraguay 2011**

Origin of virus

Argentina — 0.5

Brazil — 0.36

Bolivia — 0.14

Other — 0

0.00  0.25  0.50  0.75  1.00
Posterior probability

(D)

**Figure 4. Epidemic tracing using stochastic mapping for serotypes A and O in South America**. We show the most probable sources of serotype A epidemics in Argentina 2001 (A) and Ecuador 2002 (B). For serotype O the origins of all the Ecuadorian sequences from 2002 (C) are shown along with the origins of the strain in Paraguay 2011 (D).

(A) Serotype A



(B) Serotype O

**Figure 5. Phylogeographic generalised linear models for the spatial spread of FMDV in South America**. We show the inclusion probabilities (posterior mean) and coefficient estimates (mean and 95% credible intervals) for each of 15 predictors of spatial spread, for serotype A (panel A) and O (panel B). Vertical lines on the left subpanel of each plot show the inclusion probabilities corresponding to Bayes factors 3 and 50, while the vertical line on the right subpanel marks $\beta|\delta = 0$.

# Text S1 – Supplementary information to "Spatio-temporal Dynamics of Foot-and-Mouth Disease Virus in South America"

Luiz Max Carvalho, Nuno Rodrigues Faria, Guido König, Marc A. Suchard,
Philippe Lemey, Waldemir de Castro Silveira,
and Guy Baele

November, 2019

## Data collection and curation

### Sequence data

In Tables S3 and S4 we provide the GenBank accession numbers, date and country of collection of all sequences used in this study.

### Livestock trade data

One of the main contributions of this study lies in exploring the association between the trade of cattle, pigs and sheep and viral diffusion using a statistically sound approach that enables us to test evolutionary hypotheses concerning the influence of these predictors on viral spread (Lemey et al., 2014). In this text, we provide some additional exploratory plots of the covariate data. In Figure S3, we show the time series of livestock production among the studied countries for a time frame of over 50 years, illustrating that the production of cattle is much larger than that of any other livestock. Perhaps more importantly, while pigs and cattle production maintain a trend of positive growth throughout the second half of the twentieth century, sheep production shows a decline after 1990. Spatially, the trade of all three livestock is rather different, as shown in Figure S4. The cattle trade network is more connected, while the pig and sheep trade networks are considerably more sparse. Further, the pig network presents a clear regional pattern, with two sub-networks.

Temporally, overall trade shows remarkably different trends depending on the livestock analysed. While cattle trade shows a positive trend, pig trade is mostly stable over time and all other livestock considered show a declining trend. The FAO trade matrix is sometimes inconsistent, i.e. the exports from country $i$ to country $j$ do not always correspond to the imports into country $j$ from country $i$. In cases like these, we have picked the largest value for a given year and pair of countries $i$ and $j$.

### Vaccination and case data

We obtained the number of vaccine doses distributed in a given country in a given year starting from 1986 and ending in 2010 and divided this total by the number of live cattle produced per country and year. We could not find information for the years of 1991 and 1996 and resorted to manually imputing the missing values for these years by taking the arithmetic mean of the two adjacent years. This indicator is only a proxy for vaccination in the continent since (i) not all vaccines administered protect against all serotypes and (ii) cattle are not the only susceptible livestock. We argue our choices are justified however, because (i) most vaccine doses were actually protective against all 3 serotypes (trivalent) and (ii) cattle

1

2

are by far the most numerous livestock in South America and their vaccination is mandatory across all affected countries.

We obtained a temporal measure of vaccination by computing the average across the nine countries per year. As a measure of dispersion, we employed the standard error, multiplied by the 95% quantile of a Student t distribution with either $\nu = 8 - 1 = 7$ or $\nu = 9 - 1 = 8$ degrees of freedom (2.36 and 2.31, respectively). We show these vaccination time series in Figure S6, where it can be observed that there is substantial heterogeneity across countries in terms of temporal vaccination pattern.

# Analysis

## Maximum-likelihood phylogenies and root-to-tip divergences

In order to preliminarily assess whether the data sets used in this study contain evidence of clock-like evolution, we construct unrooted phylogenies using the maximum-likelihood procedures implemented in PhyML 3.0 (Guindon and Gascuel, 2003) under a GTR model (Tavaré, 1986) of sequence evolution with gamma-distributed site rate heterogeneity. We root the trees at the branch – as well as a position along it – that minimises the sum of squared residuals of the linear regression of root-to-tip divergences (RDV) on the sampling times. For this last step, we use the routines in the program TempEst (Rambaut et al., 2016). The resulting trend lines presented in Figure S2 are consistent with clock-like evolution, despite considerable variation between lineages, which we address later in the paper using relaxed molecular clocks.

An initial RDV analysis indicated that some sequences collected for serotype O had evolutionary divergences incompatible with their dates, and these were excluded from subsequent analyses. The excluded sequences are listed in Table S1.

Table S1: **Sequences excluded from analysis.** Serotype O sequences that displayed genetic divergences inconsistent with their sampling dates were excluded from analysis. These were sequences with older dates that presented a much higher than expected number of mutations.

| Accession | Year | Country |
|-----------|------|-----------|
| AY593818 | 1958 | Brazil |
| AY593837 | 1963 | Uruguay |
| AY593820 | 1964 | Argentina |
| AY593814 | 1965 | Argentina |
| AY593821 | 1967 | Argentina |

## Bayesian model selection

We employed generalised stepping stone sampling (GSS, Baele et al. (2016)) to compute (log) marginal likelihoods for six combinations of Markov sequence evolution and molecular clock models. A constant population size coalescent tree prior was used across all combinations. We used 50 steps with 250, 000 iterations per step, with the resulting 51 power posteriors spaced according to a Beta($\alpha = 0.3, 1.0$) distribution. Table S2 shows the results of this model selection study to determine the best combination of sequence substitution and molecular clock models for each serotype.

For serotype A, the model combination with the highest (log) marginal likelihood is a GTR sequence evolution model coupled with an uncorrelated relaxed clock assuming an underlying Gamma distribution

3

(UCG), whereas for serotype O the best model combination was GTR and an uncorrelated relaxed clock assuming an underlying log-normal distribution (UCLN). The quantile inversions for the Gamma relaxed clock model are very compute-intensive and computations for this model took much longer to finish. In the interest of computational expediency for the large number of subsequent analyses we performed in this paper, we decided to use the UCLN model for serotype A as well.

Table S2: **Log marginal likelihoods for combinations of sequence evolution (Markov) and molecular clock models.** We used generalised stepping stone sampling (GSS) from (Baele et al., 2016) to compute (log) marginal likelihoods using a constant population size coalescent tree prior. We compared the HKY and GTR substitution models, in combination with the strict clock and the uncorrelated relaxed clock models assuming either an underlying log-normal distribution (UCLN) or gamma distribution (UCG). The best fitting model for each serotype is highlighted in bold.

| Serotype | Markov model | Clock model | log marginal likelihood |
|---|---|---|---|
| A | GTR | UCLN | -13649.55 |
| | GTR | UCG | **-13619.52** |
| | GTR | Strict | -13721.67 |
| | HKY | UCLN | -13668.54 |
| | HKY | UCG | -13664.53 |
| | HKY | Strict | -13760.31 |
| O | GTR | UCLN | **-9917.44** |
| | GTR | UCG | -9925.08 |
| | GTR | Strict | -10085.06 |
| | HKY | UCLN | -9959.38 |
| | HKY | UCG | -9955.12 |
| | HKY | Strict | -10102.83 |

4

## Temporal analysis – accounting for preferential sampling

Since our sample is likely to be biased temporally – as well as spatially (but see below) –, we employ a population dynamics reconstruction method that allows accommodating the tip sampling structure, as well as the inclusion of temporal covariates (Karcher et al., 2020). The model of Karcher et al. (2020) assumes sequences are collected according to a Poisson point process with intensity

$$\lambda_s(t) = \exp\left(\beta_0 + \beta_1\gamma(t) + \beta_2 f_1(t) + \beta_2 f_2(t) + \ldots + \beta_P f_P(t) + [\delta_2 f_2(t) + \ldots + \delta_P f_P(t)]\gamma(t)\right), \quad (1)$$

where $\gamma(t)$ is the log population size at time $t$, the $\boldsymbol{\beta}$ are coefficients, $\mathcal{F} = \{f_1, f_2, \ldots, f_P\}$ are functions of interest and the $\boldsymbol{\delta}$ control interaction terms, which we do not use here. Following Karcher et al. (2020), we consider four models:

1. **Model 0:** "Naive" model, in which sampling is not accounted for (Palacios et al., 2015).

2. **Model 1:** Preferential model with $\lambda_s(t) = \exp(\beta_0 + \beta_1\gamma(t))$, denoted $\{\gamma(t)\}$ (Karcher et al., 2016);

3. **Model 2:** Preferential model using $-t$ as a covariate, thus $\lambda_s(t) = \exp(\beta_0 + \beta_1\gamma(t) - \beta_2 t)$, denoted $\{\gamma(t), -t\}$;

4. **Model 3:** Preferential model using both $-t$ and $-t^2$ as covariates, thus $\lambda_s(t) = \exp(\beta_0 + \beta_1\gamma(t) - \beta_2 t - \beta_3 t^2)$, denoted $\{\gamma(t), -t, -t^2\}$;

These models were fitted using three maximum clade credibility (MCC) trees for each serotype, obtained from three independent chains. Since results were consistent across MCC trees, we shall discuss only results for one MCC for each serotype. Complete analyses can be found at `https://github.com/maxbiostat/FMDV_AMERICA/blob/master/CODE/notebooks/BNPR_analysis_covariates.ipynb`.

In order to compare models, we employ (log) marginal likelihoods to compute (log) Bayes factors. Model 0 does not use sampling data and thus is not directly comparable to the other models. Models 1 and 2 can be seen as special cases of model 3, with suitably constructed, point-mass priors on the non-shared parameters. Therefore, we compare models 1, 2 and 3 using the (log) marginal likelihood estimate provided by the underlying integrated nested Laplace approximation (INLA, Martins et al. (2013)), which is computed following Hubin and Storvik (2016).

### Prior modelling

The smoothness of the population trajectory in the Skygrid (Gill et al., 2013) and related models is controlled by a precision parameter, $\tau$, which needs to be estimated from the data and thus be given a prior. The prior usually employed is a Gamma distribution with parameters $\alpha = \beta = 0.001$. However, we feel this prior places too much mass on small precisions (which would make estimates noisier) while at the same time allowing rather extreme values, which would artificially inflate confidence. We instead employ the penalised complexity (PC) prior from Simpson et al. (2017) (pg.14), a density from the Gumbel type II family. Let $a, b > 0$ be the shape and scale hyperparameters respectively; the probability density function is then

$$\pi_2(\tau \mid a, b) = ab \cdot \tau^{-a-1} \exp\left(-b\tau^{-a}\right), \ \tau > 0. \quad (2)$$

The authors recommend $a = 1/2$ and $b$ to be chosen such that $Pr(1/\sqrt{\tau} > S) = p$, where the value $S$ and the probability $p$ are to be chosen on substantive grounds, such that $b = \ln(p)/S$. Here we will choose $S = 1$ and $p = 0.1$, which is to say that the prior is constructed such that there is a 10% probability that the standard deviation of the log-population sizes is greater than 1. This argument leads to $b = 2.302585$. This prior has been implemented in BEAST 1.10 (Suchard et al., 2018) and can be used with the Skygrid model by using the following XML syntax:

5

```xml
<gumbelPrior shape="0.5" scale="2.30">
    <parameter idref="skygrid.precision"/>
</gumbelPrior>
```

replacing the default XML block. For the **phylodyn** package (Karcher et al., 2017), the PC prior is also implemented in this fork: `https://github.com/maxbiostat/phylodyn`, and can be used by calling `BNPR(..., pc_prior = TRUE, S, p)` and `BNPR_PS(..., pc_prior = TRUE, S, p)`.

### Simulations

In order to understand how the sequence sampling structure of both serotypes (see Figure S1) might influence results, we simulated four scenarios for each data set:

A) Using the empirical sampling dates, simulate phylogenies under a constant population size with $N_e = 10$;

B) Using the empirical sampling dates, simulate phylogenies under an exponential growth model with $N_0 = 10$ and $r = 0.1$;

C) Using a regular grid of dates spanning the range of observed dates, simulate phylogenies under a constant population size with $N_e = 10$;

D) Using a regular grid of dates spanning the range of observed dates, simulate phylogenies under an exponential growth model with $N_0 = 10$ and $r = 0.1$.

To construct a regular grid, we created a sequence of $N$ dates from youngest to oldest observed date, where $N$ is the number of sequences collected for each serotype. For all scenarios, we simulated 10 phylogenies for each serotype using the R package **timeTreeSim** (`https://github.com/maxbiostat/timeTreeSim`). We fitted the four models outlined above to all simulated data sets. Since results were consistent across phylogenies, we will discuss results for a single representative simulation per scenario and serotype. All simulations and analyses can be found at `https://github.com/maxbiostat/FMDV_AMERICA/blob/master/CODE/notebooks/`.

Typical (representative) reconstructions obtained with all four models for each scenario are shown in Figures S9 and S10 for serotypes A and O, respectively. According to these figures, the results are consistent across serotypes and preferential sampling (i.e. using the empirical dates of the sequences in this study) leads to noisier reconstructions. In addition, we notice that when the true population trajectory is exponential growth, the preferential model ($\{\gamma(t)\}$) does not recover the true population trajectory, but a model that includes time as a covariate (e.g. $\{\gamma(t), -t\}$) does. These results suggest that when population size changes over time, a preferential model that includes simple functions of calendar time can capture the true behaviour and lead to narrower credibility intervals (compare "naive" and $\{\gamma(t), -t\}$ models for all scenarios and both serotypes).

In terms of model comparison, Figure S11 shows the estimated (log) marginal likelihoods for models 1, 2 and 3 for all four scenarios. Again, results were consistent across serotypes. For scenario A, model 1 ($\{\gamma(t)\}$) attains the (log) highest marginal likelihood for all replicates. For scenario B, model 2 ($\{\gamma(t), -t\}$) performs best across replicates, likely because it is consistent with the exponential growth of the population. For scenarios C and D, model 1 is again the best fit across all replicates, showing that model selection based on (log) marginal likelihoods is parsimonious in that the simplest of model (between models 1, 2 and 3) is selected in a situation with no preferential sampling.

### Spatial analysis – phylogeographic GLM

In order to study the influence of several predictors on the spatial spread of FMDV in the continent, we employed a generalised linear model (GLM) (Dudas et al., 2017; Lemey et al., 2014). Let $P$ be the

6

number of predictors – in this study, $P = 12$ – and $K$ be the number of locations – here, $K = 8$ for serotype A and $K = 9$ for serotype O. Moreover, let $\Lambda_{ij}$ be the transition rate between locations $i$ and $j$ and $\boldsymbol{X}_{ij}$ be the relevant row of design matrix. Then:

$$\log \Lambda_{ij} = \boldsymbol{X}_{ij}^T \boldsymbol{\delta}\boldsymbol{\beta} + \epsilon_i + \epsilon_j,$$
$$\epsilon_k \sim \text{Normal}(0, \sigma^2) \text{ for } k = 1, \ldots, K, \text{ with}$$
$$\sigma^2 \sim \text{Inverse-Gamma}(0.001, 0.001),$$
$$\beta_j \sim \text{Normal}(0, 16),$$

where the $\boldsymbol{\epsilon}$ are location-specific origin and destination random effects (Dudas et al., 2017). $\boldsymbol{\delta}$ is a vector containing variables $\delta_i \in \{0, 1\}$ that describe whether predictor $i$ is included ($\delta_i = 1$) or excluded ($\delta_i = 0$) from the model and $\boldsymbol{\beta}$ are the regression coefficients. Let $p_j$ be the *prior* probability that predictor $j$ is included in the model. Here we will make the simplifying assumption that $p_j = q$ for all $j$. If we let $w$ be the probability that no predictors are included in the model, we obtain $q = 1 - w^{1/P}$. We can then fit the model and compute the *posterior* inclusion probability $\hat{\delta}_j$ from which we can compute the Bayes factor in support of the $j$-th predictor:

$$\text{BF}_j = \frac{\hat{\delta}_j}{1 - \hat{\delta}_j} \Big/ \frac{p_j}{1 - p_j},$$
$$= \frac{\hat{\delta}_j(1 + w^{1/P})}{(1 - \hat{\delta}_j)(1 - w^{1/P})}.$$

In this study we choose $w = 1/2$ to enforce a parsimonious model. See Kass and Raftery (1995) for how to interpret Bayes factors.

## Further References

Hubin, A. and Storvik, G. (2016). Estimating the marginal likelihood with integrated nested Laplace approximation (inla). *arXiv preprint arXiv:1611.01450*.

Martins, T. G., Simpson, D., Lindgren, F., and Rue, H. (2013). Bayesian computing with INLA: new features. *Computational Statistics & Data Analysis*, 67:68–83.

Palacios, J. A., Wakeley, J., and Ramachandran, S. (2015). Bayesian nonparametric inference of population size changes from sequential genealogies. *Genetics*, 201(1):281–304.

7



Figure S1: **Spatio-temporal sampling pattern of the sequences used in this study.** Circle radius is proportional to number of sequences.

Figure S2: **Root-to-tip divergence against sampling time for serotypes A and O.** Maximum likelihood phylogenies were constructed with PhyML version 3.0 (Guindon and Gascuel, 2003). We root the trees such that the linear regression between the root-to-tip divergences and sampling times yields the smallest sum of squared residuals using the program TempEst (Rambaut et al., 2016). The linear trends, along with their 95% prediction intervals (shaded) are shown.

9

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



Figure S3: **Production time series for several livestock in South America.** We show log(# of heads) of live animals for pigs, sheep and cattle, goat and horses.

10

**Cattle**

**Pigs**

under 2678.43
2678.43 – 104231.61
104231.61 – 247634.85
over 247634.85

under 96.07
96.07 – 4842.24
4842.24 – 16469.49
over 16469.49

(A)

(B)

**Sheep**

under 3103.72
3103.72 – 7690.15
7690.15 – 25377.17
over 25377.17

(C)

Figure S4: **Spatial networks of livestock trade in South America.** We represent the total trade of live animals from 1986 to 2017. Arrows connect countries if there was a non-zero number of exchanges between them. Colour scale represents total exports in number of live animals. In agreement with the data presented in Figure S3, the cattle network is the most connected, indicating that not only the number of animals has increased, but also the migration of this particular host is more frequent. Note the long range trade routes of sheep between Argentina and Colombia, which are absent from the pig network.

11



Figure S5: **Overall livestock trade in South America through time.** We show log(# of heads) of live animals traded for pigs, sheep and cattle, goat and horses. Solid lines are linear trends fitted using ordinary least squares.

(A)



(B)

Figure S6: **FMD vaccination in South America through time.** We show the number of vaccine doses per cattle head per country (Panel a) and overall (Panel b). In panel b the solid line shows the average and the dashed lines show Student t empirical intervals (see text).

13



Figure S7: **FMD cases per serotype in South America through time.** We show the number of reported FMD cases per serotype.

14



Figure S8: **Epidemic tracing of FMDV in South America – extra results**. For serotype A, we show the most probable origins of the 2001 outbreaks in Brazil (A), Uruguay (B) and Bolivia (C), along with the origins of the Colombian 2008 strain (D). Panels E and F show the origins of serotype O in Venezuela 2003 and Peru 2008, respectively.

15



(A) Constant size, empirical dates



(B) Exponential growth, empirical dates



(C) Constant size, uniform dates



(D) Exponential size, uniform dates

Figure S9: **Typical reconstructions in BNPR simulation study, serotype A**. A: constant population size $N_e = 10$ with empirical dates; B: exponential population growth with $N_0 = 10$ and $r = 0.1$ with empirical dates; C: constant population size $N_e = 10$ with uniform temporal sampling; D: exponential population growth with $N_0 = 10$ and $r = 0.1$ with uniform temporal sampling. Dashed line shows the true population trajectory, and shaded are shows the 95% BCI. Vertical tiles show the four models considered (see text for details).

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

16

(A) Constant size, empirical date

(B) Exponential growth, empirical dates

(C) Constant size, uniform dates

(D) Exponential size, uniform dates

Figure S10: **Typical reconstructions in BNPR simulation study, serotype O**. A: constant population size $N_e = 10$ with empirical dates; B: exponential population growth with $N_0 = 10$ and $r = 0.1$ with empirical dates; C: constant population size $N_e = 10$ with uniform temporal sampling; D: exponential population growth with $N_0 = 10$ and $r = 0.1$ with uniform temporal sampling. Dashed line shows the true population trajectory, and shaded ares shows the 95% BCI. Vertical tiles show the four models considered (see text for details).

(A) Constant size, empirical date

(B) Exponential growth, empirical dates

(C) Constant size, uniform dates

(D) Exponential size, uniform dates

Figure S11: **Estimated (log) marginal likelihoods for the BNPR simulation study**. A: constant population size $N_e = 10$ with empirical dates; B: exponential population growth with $N_0 = 10$ and $r = 0.1$ with empirical dates; C: constant population size $N_e = 10$ with uniform temporal sampling; D: exponential population growth with $N_0 = 10$ and $r = 0.1$ with uniform temporal sampling. Points show actual data points and boxplots show the quantiles. Only results for models 1, 2 and 3 are shown because model 0 is not directly comparable (see text).

18



(A) Serotype A



(B) Serotype O

Figure S12: **Population reconstructions for serotypes A and O under different preferential sampling models.** Reconstructions for serotype A are shown in panel A, while reconstructions for serotype O are in panel B.

19

Table S3: **Accession numbers, date and country of collection for the serotype A sequences**. When only the year of collection was known, we used the 15th of July as the collection date.

| Accession | Date | Country |
|---|---|---|
| JQ082960 | 1955-07-15 | Brazil |
| AJ306222 | 1955-07-15 | Argentina |
| AJ251476 | 1955-07-15 | Brazil |
| AY593768 | 1955-07-15 | Brazil |
| JQ082955 | 1958-07-15 | Brazil |
| AY593788 | 1958-07-15 | Brazil |
| AY593753 | 1958-07-15 | Brazil |
| JQ082961 | 1959-07-15 | Argentina |
| JQ082956 | 1959-07-15 | Brazil |
| JQ082957 | 1959-07-15 | Brazil |
| AY593769 | 1959-07-15 | Argentina |
| AY593756 | 1959-07-15 | Brazil |
| AY593789 | 1961-07-15 | Argentina |
| JQ082958 | 1962-07-15 | Venezuela |
| JQ082959 | 1962-07-15 | Argentina |
| AY593767 | 1965-07-15 | Argentina |
| JQ082962 | 1966-07-15 | Argentina |
| AY593770 | 1966-07-15 | Argentina |
| JQ082963 | 1967-07-15 | Colombia |
| AY593757 | 1967-07-15 | Brazil |
| AY593771 | 1967-07-15 | Colombia |
| AY593758 | 1967-07-15 | Venezuela |
| AJ308694 | 1968-07-15 | Argentina |
| AY593801 | 1968-07-15 | Uruguay |
| JQ082964 | 1969-07-15 | Peru |
| AY593773 | 1969-07-15 | Peru |
| JQ082965 | 1970-07-15 | Venezuela |
| JQ082966 | 1970-07-15 | Brazil |
| EU553882 | 1970-07-15 | Venezuela |
| AY593775 | 1970-07-15 | Venezuela |
| AJ306220 | 1971-07-15 | Argentina |
| AJ308695 | 1971-07-15 | Argentina |
| JQ082967 | 1975-07-15 | Ecuador |
| JQ082968 | 1975-07-15 | Peru |
| AJ409219 | 1976-07-15 | Argentina |
| EU553851 | 1976-07-15 | Brazil |
| JQ082969 | 1976-07-15 | Brazil |
| JQ082970 | 1976-07-15 | Brazil |
| JQ082971 | 1976-07-15 | Brazil |
| JQ082972 | 1976-07-15 | Colombia |
| APHA76 | 1976-07-15 | Venezuela |
| AY593787 | 1977-07-15 | Brazil |
| Continued on next page | | |

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

20

Table S3 – continued from previous page

| Accession | Date | Country |
|---|---|---|
| JQ082973 | 1979-07-15 | Ecuador |
| APHA79 | 1979-07-15 | Venezuela |
| AY593803 | 1979-07-15 | Brazil |
| KF112899 | 1981-07-15 | Argentina |
| JQ082974 | 1981-07-15 | Brazil |
| AJ306219 | 1981-07-15 | Argentina |
| JQ082975 | 1984-07-15 | Brazil |
| JQ082976 | 1984-07-15 | Colombia |
| JQ082977 | 1985-07-15 | Colombia |
| AY593794 | 1985-07-15 | Colombia |
| AJ306221 | 1987-07-15 | Argentina |
| JQ082978 | 1989-07-15 | Venezuela |
| AJ308698 | 1990-07-15 | Argentina |
| AJ308696 | 1990-07-15 | Argentina |
| AJ308697 | 1990-07-15 | Argentina |
| AJ308699 | 1991-07-15 | Argentina |
| AJ308702 | 1992-07-15 | Argentina |
| AJ308700 | 1992-07-15 | Argentina |
| AJ308701 | 1992-07-15 | Argentina |
| JQ082982 | 1997-07-15 | Brazil |
| JQ082979 | 1997-07-15 | Colombia |
| JQ082980 | 1997-07-15 | Colombia |
| JQ082981 | 1997-07-15 | Colombia |
| JQ082931 | 1999-07-15 | Peru |
| JQ082915 | 2000-05-30 | Bolivia |
| AY593782 | 2000-07-15 | Argentina |
| JQ082930 | 2000-07-15 | Peru |
| JQ082906 | 2000-08-15 | Argentina |
| AM179990 | 2000-08-15 | Argentina |
| AM179989 | 2000-08-15 | Argentina |
| AM179992 | 2000-08-15 | Argentina |
| AM179988 | 2000-08-15 | Argentina |
| AM179991 | 2000-08-15 | Argentina |
| AM179993 | 2000-08-15 | Argentina |
| AM179995 | 2000-09-15 | Argentina |
| AM179996 | 2000-09-15 | Argentina |
| AM179994 | 2000-09-15 | Argentina |
| AM179998 | 2000-10-15 | Argentina |
| AM179997 | 2000-10-15 | Argentina |
| AM179999 | 2000-11-15 | Argentina |
| JQ082937 | 2000-12-09 | Venezuela |
| KX002194 | 2000-12-31 | Argentina |
| KX002195 | 2001-02-10 | Argentina |
| KX002196 | 2001-02-18 | Argentina |
| KX002204 | 2001-03-07 | Argentina |
| KX002203 | 2001-03-13 | Argentina |
| Continued on next page | | |

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

21

**Table S3 – continued from previous page**

| Accession | Date | Country |
| --- | --- | --- |
| JQ082907 | 2001-03-15 | Argentina |
| JQ082908 | 2001-03-15 | Argentina |
| JQ082909 | 2001-03-15 | Argentina |
| AM180007 | 2001-03-15 | Argentina |
| AM180000 | 2001-03-15 | Argentina |
| AM180005 | 2001-03-15 | Argentina |
| AM180008 | 2001-03-15 | Argentina |
| AM180002 | 2001-03-15 | Argentina |
| AM180006 | 2001-03-15 | Argentina |
| AM180001 | 2001-03-15 | Argentina |
| AM180003 | 2001-03-15 | Argentina |
| AM180004 | 2001-03-15 | Argentina |
| KX002178 | 2001-03-21 | Argentina |
| KX002205 | 2001-03-27 | Argentina |
| KX002193 | 2001-03-29 | Argentina |
| AM180009 | 2001-04-15 | Argentina |
| AM180014 | 2001-04-15 | Argentina |
| AM180010 | 2001-04-15 | Argentina |
| AM180012 | 2001-04-15 | Argentina |
| AM180013 | 2001-04-15 | Argentina |
| AM180015 | 2001-04-15 | Argentina |
| AM180011 | 2001-04-15 | Argentina |
| AM180016 | 2001-04-15 | Argentina |
| JQ082910 | 2001-04-15 | Uruguay |
| JQ082911 | 2001-04-15 | Uruguay |
| KX002198 | 2001-04-20 | Argentina |
| JQ082916 | 2001-05-07 | Bolivia |
| JQ082917 | 2001-05-10 | Bolivia |
| KX002200 | 2001-05-12 | Argentina |
| KX002201 | 2001-05-15 | Argentina |
| KX002202 | 2001-05-15 | Argentina |
| AM180019 | 2001-05-15 | Argentina |
| AM180017 | 2001-05-15 | Argentina |
| AM180018 | 2001-05-15 | Argentina |
| JQ082912 | 2001-05-15 | Brazil |
| JQ082913 | 2001-05-15 | Brazil |
| JQ082914 | 2001-05-15 | Brazil |
| KX002176 | 2001-05-25 | Argentina |
| KX002177 | 2001-05-28 | Argentina |
| KX002180 | 2001-06-03 | Argentina |
| JQ082918 | 2001-06-07 | Bolivia |
| JQ082932 | 2001-06-13 | Venezuela |
| AM180020 | 2001-06-15 | Argentina |
| AM180021 | 2001-06-15 | Argentina |
| KX002181 | 2001-06-16 | Argentina |
| KX002182 | 2001-06-17 | Argentina |
| Continued on next page | | |

22

**Table S3 – continued from previous page**

| Accession | Date | Country |
|-----------|------|---------|
| KX002183 | 2001-06-17 | Argentina |
| KX002184 | 2001-07-10 | Argentina |
| AY593802 | 2001-07-15 | Uruguay |
| AY593783 | 2001-07-15 | Argentina |
| AY593784 | 2001-07-15 | Argentina |
| AY593785 | 2001-07-15 | Argentina |
| AY593786 | 2001-07-15 | Argentina |
| AY593790 | 2001-07-15 | Argentina |
| KX002186 | 2001-07-20 | Argentina |
| JQ082919 | 2001-07-23 | Bolivia |
| KX002185 | 2001-08-07 | Argentina |
| JQ082920 | 2001-08-07 | Bolivia |
| KX002189 | 2001-08-08 | Argentina |
| KX002188 | 2001-08-15 | Argentina |
| KX002187 | 2001-08-21 | Argentina |
| KX002191 | 2001-10-11 | Argentina |
| AM180022 | 2001-10-15 | Argentina |
| KX002192 | 2001-11-15 | Argentina |
| JQ082933 | 2001-12-15 | Venezuela |
| AM180024 | 2002-01-15 | Argentina |
| JQ082921 | 2002-03-30 | Bolivia |
| JQ082929 | 2002-06-15 | Ecuador |
| JQ082934 | 2002-11-28 | Venezuela |
| JQ082935 | 2003-05-14 | Venezuela |
| JQ082936 | 2003-07-03 | Venezuela |
| JQ082938 | 2003-12-12 | Venezuela |
| JQ082939 | 2004-01-23 | Venezuela |
| JQ082940 | 2004-02-11 | Venezuela |
| JQ082941 | 2004-03-10 | Venezuela |
| JQ082942 | 2004-04-01 | Venezuela |
| JQ082943 | 2004-05-19 | Venezuela |
| JQ082944 | 2004-07-12 | Venezuela |
| JQ082922 | 2004-07-15 | Colombia |
| JQ082945 | 2004-08-13 | Venezuela |
| JQ082946 | 2004-08-26 | Venezuela |
| JQ082947 | 2004-09-24 | Venezuela |
| JQ082948 | 2004-11-18 | Venezuela |
| JQ082949 | 2005-02-03 | Venezuela |
| JQ082950 | 2005-04-06 | Venezuela |
| JQ082951 | 2005-04-20 | Venezuela |
| JQ082952 | 2005-04-29 | Venezuela |
| JQ082953 | 2006-06-13 | Venezuela |
| JQ082954 | 2007-01-30 | Venezuela |
| JQ082923 | 2008-06-06 | Colombia |
| JQ082924 | 2008-06-06 | Colombia |
| JQ082925 | 2008-06-06 | Colombia |
| Continued on next page | | |

23

**Table S3 – continued from previous page**

| Accession | Date | Country |
|-----------|------|---------|
| JQ082926 | 2008-06-06 | Colombia |
| JQ082927 | 2008-06-06 | Colombia |
| JQ082928 | 2008-06-06 | Colombia |
| KU234721 | 2013-04-01 | Venezuela |

Table S4: **Acession numbers, date and country of collection for the serotype O sequences**. When only the year of collection was known, we used the 15th of July as the collection date.

| Accession | Date | Country |
|-----------|------|---------|
| APHVP1OC | 1971-07-15 | Brazil |
| AY593827 | 1971-07-15 | Venezuela |
| AJ308705 | 1977-07-15 | Argentina |
| DQ789075 | 1983-07-15 | Argentina |
| AJ308706 | 1983-07-15 | Argentina |
| KJ831663 | 1990-07-15 | Bolivia |
| AJ308707 | 1990-07-15 | Argentina |
| KJ831741 | 1992-07-15 | Brazil |
| AJ308708 | 1992-07-15 | Argentina |
| KJ831730 | 1993-07-15 | Peru |
| KJ831731 | 1993-07-15 | Peru |
| AJ292205 | 1993-07-15 | Argentina |
| AJ292207 | 1993-07-15 | Argentina |
| AJ292208 | 1993-07-15 | Argentina |
| AJ292209 | 1993-07-15 | Argentina |
| KJ831736 | 1994-07-15 | Brazil |
| KJ831742 | 1994-07-15 | Brazil |
| KJ831743 | 1994-07-15 | Brazil |
| KJ831744 | 1994-07-15 | Brazil |
| KJ831745 | 1994-07-15 | Brazil |
| KJ831746 | 1994-07-15 | Brazil |
| KJ831747 | 1994-07-15 | Brazil |
| HQ695747 | 1994-07-15 | Colombia |
| HQ695748 | 1994-07-15 | Colombia |
| HQ695749 | 1994-07-15 | Colombia |
| HQ695750 | 1994-07-15 | Colombia |
| HQ695751 | 1994-07-15 | Colombia |
| HQ695752 | 1994-07-15 | Colombia |
| HQ695753 | 1994-07-15 | Colombia |
| HQ695754 | 1994-07-15 | Colombia |
| HQ695755 | 1994-07-15 | Colombia |
| HQ695756 | 1994-07-15 | Colombia |
| HQ695757 | 1994-07-15 | Colombia |
| HQ695758 | 1994-07-15 | Colombia |
| Continued on next page | | |

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

24

**Table S4 – continued from previous page**

| Acession | Date | Country |
|----------|------|---------|
| HQ695759 | 1994-07-15 | Colombia |
| HQ695760 | 1994-07-15 | Colombia |
| HQ695761 | 1994-07-15 | Colombia |
| HQ695762 | 1994-07-15 | Colombia |
| KJ831672 | 1994-07-15 | Ecuador |
| KJ831673 | 1994-07-15 | Ecuador |
| KJ831732 | 1994-07-15 | Peru |
| KJ831733 | 1994-07-15 | Peru |
| KJ831734 | 1994-07-15 | Peru |
| AJ292206 | 1994-07-15 | Argentina |
| AJ306212 | 1994-07-15 | Argentina |
| KJ831748 | 1995-07-15 | Brazil |
| HQ695763 | 1995-07-15 | Colombia |
| HQ695764 | 1995-07-15 | Colombia |
| HQ695765 | 1995-07-15 | Colombia |
| HQ695766 | 1995-07-15 | Colombia |
| HQ695767 | 1995-07-15 | Colombia |
| HQ695768 | 1995-07-15 | Colombia |
| HQ695769 | 1995-07-15 | Colombia |
| HQ695770 | 1998-07-15 | Colombia |
| HQ695771 | 1998-07-15 | Colombia |
| HQ695772 | 1998-07-15 | Colombia |
| HQ695773 | 1998-07-15 | Colombia |
| DQ834704 | 1998-07-15 | Brazil |
| HQ695774 | 1999-07-15 | Colombia |
| HQ695737 | 2000-03-20 | Bolivia |
| HQ695744 | 2000-05-07 | Bolivia |
| HQ695806 | 2000-07-09 | Ecuador |
| KJ831738 | 2000-07-15 | Argentina |
| KJ831739 | 2000-07-15 | Argentina |
| HQ695775 | 2000-07-15 | Colombia |
| HQ695776 | 2000-07-15 | Colombia |
| HQ695777 | 2000-07-15 | Colombia |
| HQ695778 | 2000-07-15 | Colombia |
| HQ695779 | 2000-07-15 | Colombia |
| AM180025 | 2000-07-15 | Argentina |
| DQ834708 | 2000-07-15 | Bolivia |
| DQ834706 | 2000-07-15 | Brazil |
| DQ834705 | 2000-07-15 | Argentina |
| DQ834707 | 2000-07-15 | Uruguay |
| AM180026 | 2000-08-15 | Argentina |
| AM180027 | 2000-09-15 | Argentina |
| AM180028 | 2000-09-15 | Argentina |
| AM180029 | 2000-10-15 | Argentina |
| HQ695738 | 2000-11-16 | Bolivia |
| AM180030 | 2000-12-15 | Argentina |
| Continued on next page | | |

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

25

**Table S4 – continued from previous page**

| Acession | Date | Country |
|----------|------|---------|
| HQ695739 | 2001-02-19 | Bolivia |
| HQ695740 | 2001-05-05 | Bolivia |
| JN005916 | 2001-06-08 | Ecuador |
| HQ695741 | 2001-07-01 | Bolivia |
| DQ834709 | 2001-07-15 | Bolivia |
| HQ695742 | 2002-04-07 | Bolivia |
| HQ695743 | 2002-04-12 | Bolivia |
| HQ695783 | 2002-06-15 | Ecuador |
| HQ695784 | 2002-06-15 | Ecuador |
| HQ695785 | 2002-06-15 | Ecuador |
| HQ695780 | 2002-07-15 | Colombia |
| KJ831728 | 2002-07-15 | Paraguay |
| KJ831729 | 2002-07-15 | Paraguay |
| DQ834710 | 2002-07-15 | Paraguay |
| HQ695792 | 2003-01-12 | Ecuador |
| HQ695786 | 2003-01-27 | Ecuador |
| HQ695793 | 2003-02-12 | Ecuador |
| HQ695787 | 2003-06-02 | Ecuador |
| HQ695790 | 2003-06-11 | Ecuador |
| HQ695745 | 2003-07-13 | Bolivia |
| DQ834712 | 2003-07-15 | Bolivia |
| DQ834713 | 2003-07-15 | Bolivia |
| DQ834711 | 2003-07-15 | Paraguay |
| HQ695788 | 2003-07-30 | Ecuador |
| HQ695794 | 2003-08-12 | Ecuador |
| HQ695845 | 2003-09-01 | Venezuela |
| HQ695789 | 2003-09-24 | Ecuador |
| HQ695791 | 2003-11-20 | Ecuador |
| HQ695795 | 2004-01-27 | Ecuador |
| HQ695846 | 2004-05-03 | Venezuela |
| HQ695796 | 2004-06-22 | Ecuador |
| HQ695797 | 2004-06-24 | Ecuador |
| HQ695798 | 2004-06-30 | Ecuador |
| HQ695801 | 2004-07-01 | Ecuador |
| HQ695800 | 2004-07-02 | Ecuador |
| HQ695799 | 2004-07-03 | Ecuador |
| HQ695802 | 2004-07-04 | Ecuador |
| HQ695803 | 2004-07-05 | Ecuador |
| HQ695804 | 2004-07-07 | Ecuador |
| HQ695805 | 2004-07-07 | Ecuador |
| HQ695807 | 2004-07-13 | Ecuador |
| HQ695810 | 2004-07-15 | Ecuador |
| HQ695808 | 2004-07-16 | Ecuador |
| HQ695809 | 2004-07-20 | Ecuador |
| HQ695811 | 2004-07-25 | Ecuador |
| HQ695812 | 2004-07-26 | Ecuador |
| Continued on next page | | |

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

26

**Table S4 – continued from previous page**

| Acession | Date | Country |
|----------|------|---------|
| HQ695813 | 2004-07-31 | Ecuador |
| HQ695814 | 2004-09-16 | Ecuador |
| HQ695815 | 2004-09-17 | Ecuador |
| HQ695816 | 2004-09-20 | Ecuador |
| HQ695817 | 2004-10-14 | Ecuador |
| HQ695818 | 2005-01-05 | Ecuador |
| HQ695819 | 2005-03-23 | Ecuador |
| HQ695820 | 2005-04-01 | Ecuador |
| HQ695823 | 2005-04-15 | Ecuador |
| HQ695821 | 2005-04-21 | Ecuador |
| HQ695822 | 2005-04-24 | Ecuador |
| HQ695847 | 2005-05-12 | Venezuela |
| HQ695825 | 2005-05-15 | Ecuador |
| HQ695824 | 2005-05-19 | Ecuador |
| HQ695826 | 2005-05-25 | Ecuador |
| HQ695848 | 2005-06-10 | Ecuador |
| HQ695827 | 2005-06-11 | Ecuador |
| HQ695828 | 2005-06-19 | Ecuador |
| HQ695830 | 2005-06-21 | Ecuador |
| HQ695831 | 2005-06-22 | Ecuador |
| HQ695832 | 2005-06-23 | Ecuador |
| HQ695833 | 2005-06-29 | Ecuador |
| HQ695834 | 2005-07-14 | Ecuador |
| HQ695835 | 2005-07-15 | Ecuador |
| DQ834714 | 2005-07-15 | Brazil |
| DQ834715 | 2005-07-15 | Brazil |
| DQ834716 | 2005-07-15 | Brazil |
| DQ834717 | 2005-07-15 | Brazil |
| DQ834718 | 2005-07-15 | Brazil |
| DQ834719 | 2005-07-15 | Brazil |
| DQ834720 | 2005-07-15 | Brazil |
| DQ834721 | 2005-07-15 | Brazil |
| DQ834722 | 2005-07-15 | Brazil |
| DQ834723 | 2005-07-15 | Brazil |
| DQ834724 | 2005-07-15 | Brazil |
| DQ834725 | 2005-07-15 | Brazil |
| DQ834726 | 2005-07-15 | Brazil |
| HQ695829 | 2005-08-23 | Ecuador |
| HQ695836 | 2005-08-23 | Ecuador |
| HQ695837 | 2006-03-06 | Ecuador |
| HQ695838 | 2006-03-06 | Ecuador |
| DQ834727 | 2006-07-15 | Argentina |
| HQ695849 | 2006-11-21 | Venezuela |
| HQ695839 | 2006-11-24 | Ecuador |
| HQ695746 | 2007-01-23 | Bolivia |
| HQ695850 | 2007-03-19 | Venezuela |
| Continued on next page | | |

27

**Table S4 – continued from previous page**

| Acession | Date | Country |
|----------|------------|----------|
| HQ695840 | 2007-04-30 | Ecuador |
| HQ695841 | 2007-07-15 | Ecuador |
| HQ695781 | 2008-04-01 | Colombia |
| HQ695782 | 2008-04-01 | Colombia |
| HQ695842 | 2008-05-16 | Ecuador |
| HQ695843 | 2008-05-31 | Ecuador |
| HQ695844 | 2008-06-18 | Peru |
| JN005890 | 2009-03-04 | Ecuador |
| JN005891 | 2009-04-27 | Ecuador |
| JN005894 | 2009-06-04 | Ecuador |
| JN005895 | 2009-06-05 | Ecuador |
| JN005899 | 2009-06-05 | Ecuador |
| JN005896 | 2009-06-06 | Ecuador |
| JN005897 | 2009-06-06 | Ecuador |
| JN005898 | 2009-06-09 | Ecuador |
| JN005900 | 2009-06-10 | Ecuador |
| JN005901 | 2009-06-15 | Ecuador |
| JN005902 | 2009-06-18 | Ecuador |
| JN005893 | 2009-06-24 | Ecuador |
| JN005892 | 2009-06-25 | Ecuador |
| JN005903 | 2009-06-26 | Ecuador |
| JN005904 | 2009-06-26 | Ecuador |
| JN005905 | 2009-06-30 | Ecuador |
| JN005906 | 2009-07-06 | Ecuador |
| JN005907 | 2009-07-14 | Ecuador |
| JN005908 | 2009-07-31 | Ecuador |
| JN005909 | 2010-03-27 | Ecuador |
| JN005910 | 2010-04-29 | Ecuador |
| JN005911 | 2010-05-27 | Ecuador |
| JN005912 | 2010-05-30 | Ecuador |
| JN005913 | 2010-05-31 | Ecuador |
| JN005917 | 2010-06-01 | Ecuador |
| JN005914 | 2010-06-03 | Ecuador |
| JN005915 | 2010-06-03 | Ecuador |
| KX353623 | 2010-06-16 | Ecuador |
| JN005918 | 2010-06-17 | Ecuador |
| KC519630 | 2010-07-15 | Ecuador |
| JX514427 | 2011-09-15 | Paraguay |

**FGV EMAp**

## Getúlio Vargas Foundation (FGV)

Dr. Luiz Max de Carvalho
*School of Applied Mathematics*
*Praia de Botafogo, 190*
*Rio de Janeiro, RJ, 22250-900*
*Email: lmax.fgv@gmail.com*
*Phone: +55 21 3799-2348*

May 4, 2020

Dr. Santiago F. Elena
Editor-in-Chief
Virus Evolution

Dear Dr. Elena,

I would like to submit the manuscript entitled "Spatio-temporal Dynamics of Foot-and-Mouth Disease Virus in South America" for consideration for publication in *Virus Evolution*.

This manuscript was originally submitted to *Virus Evolution* in 2015 and was assigned manuscript ID VEVOLU-2015-010. After the first round of revisions, I was unfortunately unable to perform the extensive modifications requested by the referees. After concluding my PhD in 2019 I was able to return to this project and re-collect and re-analyse the data using a better analytical framework that accounts for collection (sampling) bias explicitly, including the techniques in Karcher et al. (2020) .

It is our understanding that the manuscript brings a substantial methodological advance that sheds light into the dynamics of an important livestock virus at a continental level – previous analyses have been restricted to single countries (e.g. Ecuador and Argentina) or regions (e.g. Andes). Thus, we consider the manuscript to be of interest to the readership of *Virus Evolution* both from a methodological and an applied point of view.

As an appendix to this letter I provide a point-by-point response to each point by each reviewer clarifying the improvements made. In providing our revision, we have carefully considered the helpful suggestions and critiques of yourself and three Reviewers. You will find a point-by-point response (bold) to all comments (normal text) we received. Significant changes to the manuscript find themselves in quotes.

Sincerely,

Dr. Luiz Max de Carvalho, PhD

**Editor-in-Chief**

Dear Mr. Carvalho,

Manuscript ID VEVOLU-2015-010 entitled "Spatio-temporal Dynamics of Foot-and-Mouth Disease Virus in South America" which you submitted to Virus Evolution, has been reviewed. The comments of the reviewer(s) are included at the foot of this letter.

The three reviewers have very diverse opinions, from rejection (2nd reviewer) to minor revisions (3rd reviewer), thus the decision cannot be other but to request you to perform a major revision. Therefore, I invite you to respond to the reviewer(s)' comments and revise your manuscript. Please notice that given the nature of these comments and of the required amount of modifications, the new version will be send out for a second round of review, most likely to the same reviewers.

Once again, thank you for submitting your manuscript to Virus Evolution and I look forward to receiving your revision.

Sincerely, Dr. Santiago Elena Editor-in-Chief Virus Evolution

▷**We thank the Editor-in-Chief and agree that the Reviewers' comments have helped us improve our manuscript. As the referees shared some of the same concerns, we first address those in general comments.** ◁

**General comment #1: New data**

▷ **A major criticism of our original submission was that we did not use all of the publicly available data. To address this we conducted a thorough search of GenBank, now described in the Methods section of the paper:**

We retrieved all FMDV nucleotide sequences available from GenBank (Benson et al., 2013) from the National Center for Biotechnology Information (NCBI, http://www.ncbi.nlm.nih.gov/) with more than 600 bp. This first step yielded 6,907 sequences which were then filtered to exclude all sequences that did not include the 1D (VP1) gene, resulting in 4,507 sequences being kept. We then filtered for sequences from serotypes A and O, yielding 1051 and 2350 sequences, respectively. Next, we excluded sequences that had been extensively passaged in cell culture and selected all sequences from South America (Argentina, Bolivia, Brazil, Colombia, Ecuador, Paraguay, Peru, Uruguay, Venezuela) for which information on country and year of isolation was available.

**This procedure lead to 53 additional sequences for serotype A and 43 sequences for serotype O, compared to our original submission.** ◁

**General comment #2: Accommodating temporal and spatial sampling bias**

▷**Even with broader sampling, the use of observational data brings with it the concern that temporal and spatial sampling bias might lead to incorrect inferences. We address this important concern by employing analytical methods that explicitly accommodate the possibility of sampling bias.**

On the temporal side, we employ the modelling framework of Karcher et al. (2020) to fit various coalescent-based models that specify the dependence of the sampling process on the population size or other temporal factors, while also accounting for phylogenetic uncertainty. Using marginal likelihood estimation, we can compare these models to one another and infer not only whether there is significant temporal bias but also possible explanatory factors.

To account for spatial bias, we employed a general linear model (GLM) (Lemey et al., 2014; Dudas et al., 2017) that allows for several predictors of spatial spread to be included simultaneously. This has the benefit of allowing one to construct predictors that account for possible sampling bias, such as the difference in numbers of sequences between locations and also the numbers of sequences at both origin and destination. By assessing the posterior inclusion probability of these sampling bias proxies, we can assert whether they contribute significantly to estimated dispersal rates. This framework also allows consideration of further predictors *in addition to* the sampling bias "controls".

In summary, we have updated the statistical methods in the paper so as to accommodate and test for sampling bias. It should be noted, however, that these methods are not a silver bullet; a biased sample will always be a biased sample and inferences will be affected regardless. The use of principled statistical methods helps mitigate the bias and uncover the true patterns in the data. ◁

**Reviewer #1**

The manuscript by Carvalho et al., describes the analysis of foot-and-mouth disease viruses (FMDVs) of two different serotypes (O and A) based on partial genome sequences from samples collected during a 55 year period for serotype A and 16 years for serotype O. The sequence analysis data is linked to studies on the trade of FMDV susceptible animals, e.g. cattle and pigs. The study has some interest but from my point of view, there seem to be some surprising omissions of information which may, or may not, affect the conclusions that can be drawn.

Specific points

1) In the Abstract, the text indicates that serotype O emerged (in South America) in around 1990. This is an odd statement to make since there are well known FMDV serotype O strains that predate this, e.g. O11 Campos (from Brazil in 1958), O1 Argentina (c. 1965) and O/M11/MEX from Mexico in 1952 which have all been sequenced in part and accession numbers are available (cited in Wright et al., (2013) Infect Genetics Evolution 20, 230-238). How does consideration of such sequences influence the information about the date of introduction of the viruses and their circulation in South America? In the Introduction, the authors indicate that: "Historically, serotype O has been the most prevalent serotype on the continent"? (lines 53-55, P1). The authors need to explain why the earlier strains of FMD virus were not included in their analyses.

▷We thank the reviewer for catching this. Firstly, the Mexico sample(s) was not included because we chose to restrict attention to South America. Secondly, we updated our

data sets to include many more sequences. Please see General Comment #1 for more information. ◄

2) It is curious that there is no mention of serotype C FMDV in South America.

▷Serotype C was indeed mentioned in the introduction (line 55): "Serotype C on the other hand was last encountered in the continent in 1995 in Brazil" . ◄

3) The sequence analysis is based on the VP1 coding region, this only represents about 630 nt out of a complete FMDV genome of about 8400nt. This information is not presented within the Introduction or Results section of the manuscript and only becomes apparent from Figure legends and Material & Methods (i.e. it is quite well hidden). The resolution of the analyses based on VP1 coding region sequence information alone is necessarily limited. Full genome (or near full genome) sequencing can give much higher resolution, e.g. to the level of farm-to-farm spread (e.g. see Valdazo-Gonzalez et al., (2012) PLosOne 7(11) e49650).

▷The fact that we use VP1 sequences is explicit in the methodology and all of our data and code are publicly available. We chose VP1 because it is routinely used in molecular epidemiology studies of FMDV, and by far the most abundant gene in terms of numbers of sequence on GenBank. In any case, we have now added a whole section at the end of the paper called "Limitations of this study" to address concerns of partial versus full genomes, which we agree is an important caveat. We have also included a citation of Valdazo-Gonzalez et al. (2012) in our discussion of full genomes *versus* partial sequences. ◄

4) The nature of the samples used for the virus sequence determination is also not indicated, have some been extensively passaged in cell culture? This information should be available using the accession numbers of the published sequences and clearly can influence the outcome of sequence comparisons.

▷We excluded sequences coming from samples that had been extensively passaged. The Supplementary Material includes a list of all sequences used, as well as those excluded based on this and other criteria. ◄

5) I am not familiar with the "root-to-tip" plots shown in Figure S1 but it seems to me that the slope of the line for serotype A over the period from 1995 to 2010 is not very different from that of serotype O over this limited time period and is rather different from that for the whole time period for serotype A. It would be useful if the authors commented.

▷The slope in a root-to-tip regression is a rough estimate of the evolutionary rate, and in this case both serotypes have markedly different evolutionary rates. At any rate, since the data sets have substantially changed, so have the plots. ◄

6) A major concern about identifying origins of samples is having adequate coverage of samples from potential sources. It is not entirely clear to me that the coverage of FMDV strains circulating in South America is sufficient to be able to draw good conclusions. It may be that it is but this is not clearly demonstrated.

▷Please see General comment #1. ◄

7) On P.4, lines 48-50. The text indicates that the importance of long range migration routes seems to differ for the two serotypes but the confidence values for the two serotypes seem to overlap extensively, so is this a real difference?

▷**We thank the reviewer for catching this. We have now removed this analysis as it did not address the epidemiological question appropriately. We now limit our discussion to qualitative observations of well-supported (BF $> 3$) long-range migrations for both serotypes without attempting a quantification.** ◁

8) Recombination within the coding region for VP1 alone is rather rare and thus I am not sure the check for recombination was very justified or useful (P. 9 lines 8-11).

▷**We have removed this analysis.** ◁

**Reviewer #2**

Comments to the Author The "Spatio-temporal dynamics of FMDV in South America" study by Carvalho et al. describes the historical spatio-temporal dispersal of FMDV across the South America continent reconstructed using phylogeography analyses employed through a Bayesian spatial diffusion model. The reported results define transmission networks and transmission hubs (at country level) that would explain the historical spread of FMD within the continent. In addition, the authors deal with variables likely associated with the FMD spread (i.e. geographical distance, livestock density, and livestock trade) in trying to explain their causative effect associated with historical FMD outbreaks. As last attempt, the authors correlated the demographic dynamics of FMDV with reported number of FMD outbreak and vaccination coverage in trying to assess the impact of FMD control policies on the FMDV diversity and its population expansion/contraction through time. The paper has been already published as an arXiv (http://arxiv.org/abs/1505.01105) the 5th of May. Although the methodological approach has been previously used in different setting and the results would be interesting for a computational basis, there are several aspects of the study that need to be carefully considered before the paper would be suitable for publication. In fact, the presence of bias in the data used is potentially producing an incorrect picture of FMD in South America. In addition, although the authors are examining the potential impact of the sampling bias in the sequence data analysed, this is only properly discussed as Supplementary Text and not in the main paper, where they assume the results as correct, valid and without bias. One major problem of the study is the data. They claim to have analysed all the data publicly available in GenBank but (as detailed below) this is not correct and the analyses should be repeated including the full sequence dataset. I would be, therefore, really cautious to draw important conclusion from this study given the issues reported, which might hold true only for the time-frame pictured from the data you have analysed. As already said, this study need a proper revision before being published and this main revision would involve the re-analysis of all the data adding all the sequences available in GenBank and which are not included in this version.

▷**We thank the reviewer for such a through assessment of our work. It seems the reviewer's main concerns were (i) incomplete sampling of available data; (ii) sampling bias (even in face of all available data) and (iii) the level of generality of the results in face of (i) and (ii). We have now re-done the data collection and improved the data**

**sets we analyse. We also employ better models that account for sampling bias both in time and space. For more elaborate information, please see tje responses below and also General Comments #1 and #2 above. ◄**

Problems of the study: - The authors claim to have used all publicly available VP1 sequences from GenBank, but after inspection this is not true. In fact, there are quite a number of sequences that has not been included in the study. If this has been done intentionally, the reason for this decision should be discussed in the paper; if not, I strongly suggest checking better in GenBank what is missing from your analyses (you have even the GenBank Accession Nos of the missing data in one of your reference [Malirat et al., 2007]). For the serotype O, for example, you are totally "ignoring" the sequences before the 1994 but, what about the O/Campos (O/Br/58) and the Argentinian samples of the 82-83 or the Caseros/67, the Selab/77? I could make the list much longer. This is valid for the serotype A as well (e.g. among others, where is the A10/Arg/61?). Therefore, if you want to claim to have analysed the complete VP1 coding sequence data for South America you need to re-perform all the analyses, because the results might provide you a completely different picture of FMD in South America (see your conclusion on the Colombian origin of the type O in South America).

**▷We thank the reviewer for their careful assessment of our data. It is indeed true we had not included many sequences that could otherwise have been analysed. This has now changed and we have collected 53 additional sequences for serotype A and 43 sequences for serotype O (See General Comment #1). ◄**

- All the results discussed on the spatial diffusion of FMDV in the whole South America continent should be treated with cautions (potentially providing you incorrect data), considering that you have: missing information from missing sequences; a sampling bias in your data according to time and country. It is well known that FMDV was introduced (as you pointed out in the introduction) by human migration from Europe in the end of the 19th Century with early reports in Argentina between 1860 and 1870, and 1895 in Brazil. There were at least two distinct introductions in the North and one in the South but, before the 1922, it is really difficult to say which serotype was (i.e. before the FMDV typing was performed). However, it is clear that the early spread of FMD was coming from the South. Argentina at the time was one of the main export hubs of livestock to the continent and even the Mexico outbreak in 1922 has been attributed by the introduction of infected cattle from Argentina. This holds true for: Chile, 1920 outbreak (decline of cattle industry during 1912 in Chile with large introduction from Argentina); official report in Venezuela 1950 (potential from importation of Argentinian meats/livestock back to the 1947). It might worth to know that the Andes acted as a barrier for taking FMD out of Chile and the western part of South America, until when the regional animal movements and trade primarily caused the spread. The countries of the Rio de la Plata, which were sharing the Pampas ecosystem, experienced an early wave of disease spread and by the 1920 FMD was in Uruguay, Paraguay and Brazil. Historical data suggest that type O was introduced most likely from the South (maybe Argentina) and type A introduced from Europe. From your analysis the initial historical FMD wave has not been characterised (in the years before 1994 for type A; very limited and potentially biased before 1965 for type O). The only part which might be more realistic is the FMD transboundary movements within the "countries triangle" of Colombia, Ecuador

and Venezuela, that could sounds more like from Argentina-Venezuela-Colombia-Ecuador, even though you have the effect of sampling bias that needs to be discussed.

▷**We share the reviewer's concerns that sampling bias might be an important factor in our analyses. This has been addressed using state-of-the-art phylodynamic methods which accommodate both temporal and spatial sampling bias. While not a panacea, these methods are the best one can do in face of incomplete and potentially biased data. For more information, please see General Comment #2.** ◁

- Although you presented some data on sampling bias in your Supplementary Text (but this might be not satisfactory enough given the problem in the dataset), there is a real problem of sampling bias and this need to be addressed in your main text as well. How does the model deal with missing links? Is the prediction robust enough to provide a clear indication of virus spread in such a large geographical range? This could be a serious limitation (and problem) of your study. For serotype A, you analysed 131 VP1 sequences of which 44% are from Argentina (of which ~70% are from 2000 and 2001) and 21% from Venezuela (of which ~86% are recent samples - after 2001). In addition, the majority of your oldest samples are only from Brazil and Argentina. For serotype O, you have 167 sequences in total of which ~54% are from Ecuador (all after the 2002 and have been previously analysed - along with 30 sequences that have been included in this manuscript as well). Therefore, you have 90+30=120 sequences already analysed in a previous paper. Among the other, 36 sequences from Colombia (~22% of the total) are barely covering the 2000s (as you claim 1994 to 2008), since you have 5 sequences from 2000, 1 from 2002 and 2 from 2008, a gap of 6 year. For the type A database, your oldest samples are only from Colombia. You attempted a random sub-sampling that, as far as I understood, have not taken into account the time of sampling, but just the quantity of data from each country. Maybe you need to account for time in your sub-sampling.

▷**We now consider a model that explicitly accounts for temporal biases in the sampling process.** ◁

- When doing analysis on sequences extracted from GenBank a detailed list of the sequences with their GenBank Accession No (along with associated metadata) should always be provided. Although a webpage (but this is difficult to check and, probably, the majority of the readers would not bother to access to your website) has been set up for the paper there are no GenBank references for the serotype A. This information should be included either as a table in the main text or as a S3.

▷**This list is now available in Supplementary Material.** ◁

- This paper has been already published as an arXiv (http://arxiv.org/abs/1505.01105) which has been submitted the 5th of May and updated the 2nd of June (I received this review the 12th of June). Although it is common for theoretical maths, physics and computer sciences studies to be published as arXiv before being properly peer-reviewed, this is not the case with study dealing with topics as in this case. Since the study needs a substantial review and re-analysis of the data, I would suggest to withdraw your submission to arXiv

▷**Our submission of the initial version of our manuscript to arXiv does not pose any conflict with the journal's policy. Additionally, arXiv and bioRxiv are well known at this point to host pre-prints that have not yet been peer-reviewed, a practice that is commonly accepted in our field of research (as can be seen from the many pre-prints circulating without peer review on SARS-CoV-2).** ◁

Major Comments: - Page 1 Line 36: "Our dating". This result is only compatible with your dataset and must not be related with the incursion of FMD in South America. Your dating is referred to the MRCAs of the data you have analysed but not the MRCAs of both the type A and O clades in South America. As already detailed the occurrences of both serotypes are much earlier than your estimates, which therefore are misleading in the description. If you comment on the South America FMD phylogenetic history, you need to do that only in line with the data you analysed and not as a general picture.

▷**Whilst the occurrence of both serotypes might have occurred much earlier than the estimates we obtain, the estimates are for the origin of the *circulating* strains. It is entirely possible FMDV has been introduced several times in South America.** ◁

- Page 2 Line 7: "By the 1970s". Again this is not true. During 1950s FMD was already causing problems in Argentina, Brazil, Chile, Peru, Uruguay, Venezuela, Colombia, and Ecuador. This picture is larger than a regional scale.

▷**"Causing problems" is not the same as having widespread epidemics. See the reference we give (Saraiva, 2003) for more details.** ◁

- Page 2 Line 33: "using all". You are not using all the sequences available in GenBank. For example and as already commented, I cannot find the O Campos in your fasta file (and this is only one). I strongly suggest doing a better search and re-perform all the analyses.

▷**This has now been done.** ◁

- Page 3 Line 20: "the time of the most recent". You need to clearly state here that these estimates hold true only for your sequences analysed (MRCSs of the data) and not of the entire South America because saying that is misleading and incorrect.

▷**We have now clarified that the estimates as presented pertain to the data set at hand.** ◁

- Page 3 Line 21: "indicating a more recent origin". Again this is only true for your data and should be stated. Serotype O outbreaks have been reported in South America since the initial wave, but clear reports start from 1940-50 (e.g. massive outbreak in Peru in 1962; 1950 official report in Venezuela; 1957 A, O and C in the entire Rio Grande do Sul).

▷**Again, this is the origin of the circulating strains, which is all that can be said from any particular sequence data set.** ◁

- Page 3 Line 22: the results show a faster clock for the type O than the A (even this is not really a lot faster - considering the VP1 only there is a difference between the two of ~4nt changes per year). Might this be due to the different molecular clock model used for type A and O?

▷**The difference is not due to the choice of model because we considered the same set of molecular clock models for both serotypes, selecting those that provided the best fit for each. Rate estimates are, however, widely consistent across models, for each serotype.** ◁

- Page 4 Line 49: "Remarkably". You claim that there is a difference between long-range migration routes between serotypes (reported as 0.14 and 0.05 for type A and O, respectively). However, the 95% intervals are really similar and both containing the zero, I should then say that this is not so remarkable.

▷**This analysis has now been excluded from the manuscript.** ◁

- Page 4 Line 57: "The most probable". This result might indicate that the 2001 reappearance of FMDV in Argentina was a persistent virus foci (maybe carrier?) or maybe some missing links (i.e. sampling bias) exists in your data which are not including contemporary isolates from neighbouring countries (besides Brazil) and, therefore, this would impact on your results. Since this would be quite an interest topic (even though only on a retrospective line), you need to discuss this in more details, assessing as well the validity of your results.

▷**Our GLM analyses attempt to account for sampling bias, and find only mild evidence of bias (BFs < 3 indicate weak support, see Figure 5 in the revised manuscript). We have now expanded the discussion of these issues a little more in the revised manuscript.** ◁

- Page 5 results on the Venezuelan origin of Andean FMDV spread: Considering that you have a bias in your sequences for type A, you need to really consider with caution your results and discuss more about the impact it might cause.

▷**The question of bias has now been addressed more thoroughly (see General Comment #2).** ◁

- Page 5 Line 16: "Similar to what was found for Venezuela". You present data on Colombia saying that this results is similar to type A for Venezuela? however, I would rather imagine that you are discussing about source of type O in the north from Colombia. Is this true? If so, please rewrite the sentence to make that clear. In addition, since you have all the oldest historical samples of type O from Colombia, the logical reasoning would be that of course the analysis point to Colombia as the main transmission hub. But, what about the sequences you have not included in the analysis? Should this provide you a different picture? You need to clearly discuss this issue (and of course re-perform the analysed including all the data available)

▷**The discussion of the spatial origins results has now been completely re-written, among other reasons to accommodate the new sequences being analysed. While common sense would dictate that the location with the oldest samples would be inferred as root, this is simply not true for the CTMC model we employ: the root state can potentially be any of the sampled states (countries).** ◁

- Page 5 Lines 34-38 and following paragraph (Lines 41-57): "or serotype A". This sentence is really confusing and need to be better formulated. For type A you find that geographical distance drives the diffusion, whilst this has a higher statistical support for type O but not

like the cattle exchange. Now, the question is, what cattle exchange means? This implies geographic distance as well (because you are defining trade between countries, which in South America are not so very close), isn't it? So, the geographic distance is the main effect of FMDV diffusion or a confounding effect for cattle trade? I am really struggling to find a logic behind this results (or its analytical approach) considering that you have a strong bias in your data (both spatially and temporally) and you are analysing the geographical distance and trade (both cattle and swine) variables separately? have you checked for multicollinearity?

▷**The confusing sentence has now been re-written. We now employ a GLM approach to modelling the factors associated with spread**, **which allows us to account for multicollinearity and sampling bias simultaneously.** ◁

- Page 6 Line 9: Sensitivity analysis. The sub-samples (as referred in table S4 and S5) is excluding the over-represented countries, i.e. Argentina and Colombia. However, if you exclude Argentina from the type A data, you have now Venezuela that is over-represented (the same holds true for type O, for which Colombia is the over-represented after the exclusion of Ecuador). Since both the Argentinian and Ecuador samples are, let's say, monophyletic (collected for the majority within epidemics), this are not really changing the global picture. In addition, you claim (Page 24 Line 37) that removing Argentina from the type A analysis move the MRCA estimate of ~6 years. Is this because you eliminate one of the oldest sequences present in your data, thus leaving only the Brazil '58 and, therefore, introduce a more substantial sampling bias/uncertainty? For type O, considering that the oldest samples are from Colombia, removing Ecuador has no impact in the results. I am getting confused to understand which methodology is behind your random sampling approach used for the 5 sub-sampling. Is this a proportional random sampling (with a temporal sub-sampling as well) of each country?

▷**The sub-sampling analyses have now been removed from the manuscript. The reasons for this are two-fold: first, subsampling has the unfortunate property of exploding combinatorially in the number of strata one wants to consider. Secondly, it would be hard to concatenate results from several hundreds of replicates in order to assess whether sampling bias had a role. Instead**, **we now adopt a modelling framework that accounts for temporal and spatial sampling bias directly, in a model-based fashion.** ◁

- Page 6 Demographic reconstruction: You are discussing the increase/decrease in FMDV diversity according to the reported activity of FMD and the control policies (i.e. vaccination) imposed. However, it seems that it is difficult to correlated like-with-like in your graph(s): you have doses/head of vaccine (this could be monovalent, bi-, tri-; strain(s) used), no of FMD cases (I suppose reported no of outbreak - this could be 1 individual of 1000s of animals infected) and viral diversity. One point that is completely missed in your discussion is the vaccine efficacy and this might impact in your analysis (i.e. some reports of drop in efficacy of the O campos vaccine). In addition, you comment that after 2001 an increase in vaccine doses resulted in a decrease in viral diversity: although from the FMD outbreak data is true for type A, it is not clearly valid for type O, which maintained a more stable trend (of course with some fluctuations). Is this, again, an issue due to bias in your data? A previous study (de Silva et al., 2012) describes how BSP incorrectly reconstructed a decrease in the last part of a datum epidemic when the population was still growing. This problem was

related to the lack of genealogical information at later times. Would this be the case for your analysis as well?

▷**We thank the reviewer for this astute observation. The preferential sampling models considered in the revised manuscript do show that the corrected $N_e(t)$ plots (Figure 2, right panel) are somewhat different from naive estimates.** ◁

- Page 7 Line 31: "the inclusion of archival". You are discussing about the impact of using an outgroup into your analysis and fail to analyse that (although you could easily extract some sequences from GenBank). In addition, you claim that the type O was circulating in the continent with its root in Colombia but you do not include any samples prior to the 1994. This analysis is incorrect and should be appropriately revised.

▷**This issue has already been addressed in previous responses.** ◁

- Page 7 Line 55: "Previous studies". It seems that the study you referred describes similar cycle of 4-5 years for both serotypes and, moreover, this would be really complicated and dangerous to apply as a general rule (since it is a country-based estimate). I would suggest deleting this sentence. In addition, from you skyride plot it is difficult to say that a 4-5 year FMDV cycles exists.

▷**Done.** ◁

- Page 8 Line 8: "The diversity bottleneck". You commented about the bottleneck for the type A diversity reconstruction as an effect of FMD epidemics affecting several countries after the 2000, but I would remind you that the majority of your samples (collected mainly from epidemics) are from the 2000 afterwards. Therefore, this again would be a confounding effect due to bias in your data (i.e. is the skyride estimate affected by the number of coalescent events in your phylogeny?).

▷**The new preferential sampling models should accommodate this (see Figure 2 in the revised manuscript).** ◁

- Page 8 Lines 18-21: "Our results suggest". This is incorrect and needs to be properly assessed when a more comprehensive analysis, which would include all the type O isolates, has been performed. If the paper of Carvalho et al., 2013, indicates the same results (i.e. describing the origin of the FMDV serotype O in South America from Colombia using the very same data), that needs a proper review as well.

▷**It is not incorrect to say that our results suggest a particular inference. Ultimately, there is a limit to what can be said from limited observational data; we do our best to accommodate potential biases in the sequence sampling.** ◁

- Page 8 Line 26: "viral effective size". This is a reminder about the previous comment on Demographic reconstruction.

▷**Acknowledged.** ◁

- I am not familiar with the methodology behind that, but I suppose that the analysis of epidemiological predictors (i.e. cattle and pig trade/livestock data) seems to have been constructed around an "average" value which potentially does not describe the space-time trends of trade routes and animal movements. Does this have an impact on your results generated? If so, what is the validity of those results? Please, comment on this.

▷**We now employ a general(ised) linear model (GLM) approach that allows us to break up trade into temporal chunks as well as include all predictors at once.** ◁

- You might consider using some sequences (e.g. O BFS 1860) as outgroup for your phylogenetic reconstruction and perform a more detailed analysis that would include all your sequences (maybe using a random local clock to account for variability in the rates), therefore shaping the entire tree topology. In addition, the phylogenetic trees, as are presented now, are confusing and really difficult to read.

▷**The analyses presented here pertain to rooted time-trees. As such, the suggestion of using an outgroup does not apply.** ◁

Minor Comments: - Page 1 Line 34: "environmental". Are you really using environmental data (e.g. air and bathing water quality)? Or do you mean livestock population and trade data, so more epidemiologically-related data or population data?

▷**We thank the reviewer for this comment. We now refer to the data collected for this paper as epidemiological and populational, excluding "environmental".** ◁

- Page 1 Line 42: "Our findings". This is a general sentence and might lead to the assumption that evolutionary and spatial dynamics of serotype A and O are globally different (which might be not the case). Just highlight the South America setting.

▷**The sentence has been re-written.** ◁

- Page 1 Line 50: ", the most important". Is FMD the most important animal disease or, better, is one of the most?

▷**It seems to be the case, specially for countries such as Brazil and Uruguay which depend on meat exports.** ◁

- Page 2 Line 17: References 4 and 14 are duplicates.

▷**We thank the reviewer for catching this.** ◁

- Page 2 Line 17: Use Di Nardo et al. [12].

▷**Done.** ◁

- Page 2 Line 18 and Line 20: "environmental". Check the meaning of "environmental data" with what you are trying to analyse.

▷**Done.** ◁

- Page 2 Line 28: "in the continent". Which one? Do you mean at "continental" level? Or you are only referring to South America?

▷**We are referring to South America, which is a continent so both statements would be correct.** ◁

- Page 3 Line 48: "..we employ an asymmetric". You have already detailed your analysis procedures in the Material and Methods section. This could be deleted.

▷**Done.** ◁

- Page 5 Line 12: "We provide evidence of Venezuela.." Which region you are referring to? The Andean region? Please, specify.

▷**Addressed.** ◁

- Page 5 Line 26: "trade and viral diffusion, we collected". I think it would be better to say "we obtained data from" because maybe you have not been in the field collecting data.

▷**Corrected.** ◁

- Page 7 Line 22: I would rather use: "see Figure 3 in [6]" or "see FMD historical outbreak data in [6]".

▷**Re-written.** ◁

- Page 8 Line 7: "..for viral Ne in both". You previously discuss about viral diversity and now present the effective population size (Ne). Since the skyline family is based on the ?=Net, you should describe Ne only if you have extracted that estimate with the appropriate measure of generation time, which you are not having or even discussing (i.e. in your graph you should report in your y-axis legend the compound value as well - Net). I suggest using viral diversity throughout the paper.

▷**This has now been resolved. Thank you.** ◁

- Page 9 Line 13: It seems you used BEAST 1.7.5 - or even and older 1.7.2 version - (from your web available .xml) to perform the analyses, although a more recent version is available (1.8.2). Is the latest version more robust and efficient in the results generated? Have you re-analysed your data using the latest version and producing similar results?

▷**All analyses have now been conducted with the latest stable version of BEAST (v1.10).** ◁

- You have a type O sequence from Peru (GenBank Accession No. HQ695844.1) you say collected in 2004, whilst in your previous paper on Ecuador is defined as 1994. I checked in GenBank and this is from 2004. Therefore, you need to amend your previous paper with the correct date.

▷**We thank the reviewer for catching this.** ◁

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

**Reviewer #3**

Comments to the Author: Your reviewer is not an expert in phylogeographic or phylodynamic analyses, but does have experience in Bayesian analyses in general and FMDV epidemiology and evolution.

The authors present an analysis of the spatio-temporal dynamics of FMDV in South America. They conclude that serotypes O and A behave quite differently, with different rates of evolution, different circulation networks, and different predictors of spread. Their analysis is based on ~300 VP1 sequences retrieved from public databases, and combines a series of Bayesian phylogenetic/geographic/dynamic analyses to come to its conclusions.

It is unfortunate that there is such limited sequence data to work with, especially since it is only for VP1, a very small component of the FMDV genome (albeit a very variable one), but I am concerned more by other problems with the underlying data on which it is based, and the soundness of the resulting conclusions. Assuming these concerns can be addressed, however, the paper provides new conclusions on the spread of the disease in South America, will be of interest and value to the FMD research community, and should be published.

▷**We thank the reviewer for their assessment. Some of the concerns raised have already been addressed in previous responses.** ◁

This would be especially true is there were any evidence in terms of the known epidemiology of the disease which might support the assertion that it differs so strikingly between the serotypes, such as a differential likelihood in different serotypes of airborne transmission (which might favour geographical spread?) or subclinical infection and subsequent transmission (which might favour transmission despite inspection as a result of trade?).

▷**We thank the reviewer for this important, thought-provoking comment. An alternative explanation to the differences observed might be the stochasticity involved in the epidemics: even minor differences in transmissibility or incubation period, say, might be amplified in terms of attack ratio and other population-level variables. In other words, it might be the case that highly variable population processes are actually responsible for most of the observed variation. We have now added a paragraph at the end of the Discussion expounding a bit more on this topic. We again thank the reviewer for reminding us of this important caveat.** ◁

My specific major concern is with the effect of differential sampling effort resulting in different numbers of samples in each country, rather than this being a feature of the epidemiology of the disease. Chile, Guyana, French Guiana and Suriname have no FMDV sequences, and Paraguay has no A sequence even though wrlfmd.org shows that Paraguay, Guyana and Chile have experienced recorded outbreaks. It is very likely that French Guiana and Suriname have too. Peru, Paraguay and Uruguay also have very low numbers of samples in total. Some work has been carried out to investigate the sensitivity to spatial sampling heterogeneity, but not with respect to the predictors of spread. Depending on whether there is detailed information on the location of VP1 sequences within countries, this inference could depend strongly on these poorly sampled (or unsampled) countries, and it would seem important to investigate whether this effect alters the conclusions, perhaps by removing

the poorly sampled countries (since we can't add in potential missing countries), and just investigating spread between the countries with high sample numbers for the serotype.

▷**This is a legitimate concern. As discussed previously, we prefer a "complete-data" approach to assessing sampling bias: we fit a GLM with a stringent, sparsity-inducing prior, and include as many predictors associated with sampling bias (difference and product in number of sequences, number of sequences as origin-destination predictors) as possible. The analyses in Figure 5 in the revised manuscript shows that none of the "sampling bias" predictors achieves an appreciable level of support (Bayes factor).** ◁

It would also seem important to investigate whether the reason for the difference in observed predictors relates not to the different epidemiology of the serotypes, but to the different control policies in the different time periods studied. This could be investigated by doing all of the inference for serotype A again based just on the recent data.

▷**This again is a very astute observation, for which we thank the reviewer. We posit that the (temporal) preferential sampling models should capture this effect -- should it exist -- albeit imperfectly, through the coefficient of $-t$. Unfortunately the present state-of-the-art methods do not allow for an easy way of incorporating complex (non-simple) time-covariates.** ◁

Also, while it seems plausible to suppose that there might be a different substitution rate, this difference seems high, and a casual inspection of Figure S1 also suggests that there might be much faster substitution rate in A during the period for which O data exists, though I don't know why that might be the case.

▷**Differences this large do exist among FMDV serotypes (see for instance Tully and Fares (2008)). We agree that some of the differences could be caused by different sampling, etc, and that is why the substitution rate should not be read too much into. See Holmes et al. (2016) for a nice discussion on why differences in substitution rate may not and often do not reflect differences at the replication level.** ◁

Finally, I see no explanation or justification for why the molecular clock model should be different between two serotypes of the same virus circulating through the same species in the same region. Some explanation would seem appropriate, or an investigation of what the implications for other conclusions might be if this might not be the case.

▷**See our comment above. We have now added a reference to the Holmes et al. review and added a couple sentences to the Discussion further reinforcing this point. Thank you for drawing our attention to the issue.** ◁

Minor detail:

Kullback-Leibler is occasionally misspelt as Kullback-Liebler.

▷**Fixed. Thanks.** ◁

References

Benson DA, Cavanaugh M, Clark K, Karsch-Mizrachi I, Lipman DJ, Ostell J, Sayers EW. 2013. GenBank. *Nucleic Acids Res*. 41:36–42.

Dudas G, Carvalho LM, Bedford T, et al. (11 co-authors). 2017. Virus genomes reveal factors that spread and sustained the ebola epidemic. *Nature*. 544:309–315.

Holmes EC, Dudas G, Rambaut A, Andersen KG. 2016. The evolution of Ebola virus: Insights from the 2013–2016 epidemic. *Nature*. 538:193–200.

Karcher MD, Carvalho LM, Suchard MA, Dudas G, Minin VN. 2020. Estimating effective population size changes from preferentially sampled genetic sequences. *PLoS Computational Biology*. .

Lemey P, Rambaut A, Bedford T, et al. (11 co-authors). 2014. Unifying viral genetics and human transportation data to predict the global transmission dynamics of human influenza H3N2. *PLoS Pathog*. 10:e1003932.

Saraiva V. 2003. Epidemiology of Foot-and-mouth disease in South America. In: Dodet B, Vicari M, editors, Foot and mouth disease: control strategies, Paris: Elsevier SAS, pp. 43–54.

Tully DC, Fares MA. 2008. The tale of a modern animal plague: tracing the evolutionary history and determining the time-scale for foot and mouth disease virus. *Virology*. 382:250–256.