1   **Examining the molecular clock hypothesis for the contemporary evolution of the rabies**

2   **virus**

3   **Authors**: Rowan Durrant[1]*, Christina A. Cobbold[1,2], Kirstyn Brunker[1], Kathryn Campbell[1],

4   Jonathan Dushoff[3,4,5], Elaine A. Ferguson[1], Gurdeep Jaswant[1,6,7,8], Ahmed Lugelo[1,8,9],

5   Kennedy Lushasi[8], Lwitiko Sikana[8], and Katie Hampson[1]

6   **Affiliations**:

7   1.  Boyd Orr Centre for Population and Ecosystem Health, School of Biodiversity, One

8       Health & Veterinary Medicine, College of Medical, Veterinary & Life Sciences,

9       University of Glasgow, Glasgow, UK

10  2.  School of Mathematics and Statistics, University of Glasgow, Glasgow, UK

11  3.  Department of Biology, McMaster University, Hamilton, Ontario, Canada

12  4.  Department of Mathematics and Statistics, McMaster University, Hamilton, Ontario,

13      Canada

14  5.  M. G. DeGroote Institute for Infectious Disease Research, McMaster University,

15      Hamilton, Ontario, Canada

16  6.  University of Nairobi Institute of Tropical and Infectious Diseases (UNITID),

17      Nairobi, Kenya.

18  7.  Tanzania Industrial Research Development Organisation (TIRDO), Dar es salaam,

19      Tanzania

20  8.  Environmental Health and Ecological Sciences Department, Ifakara Health Institute,

21      Ifakara, Tanzania

22  9.  Global Animal Health Tanzania, Arusha, Tanzania

23  **\*Corresponding author**: Rowan Durrant: r.durrant.1@research.gla.ac.uk.

24

25 **Abstract**

26 The molecular clock hypothesis assumes that mutations accumulate on an organism's genome

27 at a constant rate over time, but this assumption does not always hold true. While modelling

28 approaches exist to accommodate deviations from a strict molecular clock, assumptions about

29 rate variation may not fully represent the underlying evolutionary processes. There is

30 considerable variability in rabies virus (RABV) incubation periods, ranging from days to over

31 a year, during which viral replication may be reduced. This prompts the question of whether

32 modelling RABV on a per infection generation basis might be more appropriate. We

33 investigate how variable incubation periods affect root-to-tip divergence under per-unit time

34 and per-generation models of mutation. Additionally, we assess how well these models

35 represent root-to-tip divergence in time-stamped RABV sequences. We find that at low

36 substitution rates (<1 substitution per genome per generation) divergence patterns between

37 these models are difficult to distinguish, while above this threshold differences become

38 apparent across a range of sampling rates. Using a Tanzanian RABV dataset, we calculate the

39 mean substitution rate to be 0.17 substitutions per genome per generation. At RABV's

40 substitution rate, the per-generation substitution model is unlikely to represent rabies

41 evolution substantially differently than the molecular clock model when examining

42 contemporary outbreaks; over enough generations for any divergence to accumulate, extreme

43 incubation periods average out. However, measuring substitution rates per-generation holds

44 potential in applications such as inferring transmission trees and predicting lineage

45 emergence.

46

47 **Author Summary**

48 Rabies is a neglected disease that kills around 60,000 people each year. After entering the

49 body, the incubation period of the virus is usually less than one month, but can sometimes

50 span months to years. While we normally assume a virus accumulates mutations at a constant

51 rate, it is possible that rabies' occasional long incubation periods mean that mutations

52 accumulate at varying rates if the virus replicates (and thus mutates) more slowly during the

53 incubation period. We compared how the rabies virus evolves over time using two simulation

54 models where mutations either occur per unit time or per infection generation. We also

55 calculated the mean substitution rate per infection generation, which can be useful for

56 inferring linkage between related rabies cases. We found that at realistic substitution rates for

57 the rabies virus, we could not distinguish between the two models. Our calculations show that

58 in most generations no mutations are expected to occur. Thus, over a time period long enough

59 to observe genetic divergence, occasional long incubation periods would be "cancelled out"

60 by shorter than average incubation periods, meaning that the two models are almost

61 equivalent. However our work suggests that modelling substitution rates per generation may

62 be useful for epidemiological inference.

## Introduction

64 The molecular clock hypothesis assumes that the genomes of organisms accumulate neutral

65 mutations at a constant rate over time, either across all lineages (the "strict molecular clock")

66 or within each individual lineage but with some degree of variation between them (clock

67 models with this assumption include the relaxed and multirate clock models) (1–3). The

68 ability to sample viral sequences through time, and the application of the molecular clock

69 hypothesis to these sequences, has led to massive advances in using viral genetic data to

70 investigate disease outbreaks (4). The clock rate, measured in substitutions per site per unit

71 time, can be used to estimate how long ago pathogens diverged (5), and the date of infection

72 of individual infected hosts (6). Combining the analysis of epidemiological and genetic data

73 has allowed further insights into the history of outbreaks (7), and the introduction of

74 geographic data provides estimates as to rates of spread and the frequency and source of

75 introductions (8,9). However, in order to conduct these phylogenetic analyses, genetic

76 divergence must increase appreciably over time in the dataset under investigation (10).

77 Whether or not the viral population is measurably evolving, and thus whether it contains

78 sufficient temporal signal for phylogenetic analysis, depends mainly on the evolutionary rate,

79 the sequence length and the length of time sequences are sampled over being sufficiently

80 high. Various methods exist to assess temporal signal, the most commonly used being root-to-

81 tip divergence plots (11,12) implemented in tools such as TempEst (13) , but these also

82 include Bayesian evaluation of temporal signal (BETS) (14)  and the date-randomisation test

83 (15).

84 The rabies virus (RABV) is a negative-strand RNA virus, with a genome size of

85 approximately 12 kilobases. While RNA viruses generally have high mutation rates due to a

86 lack of proofreading by RNA polymerases, RABV has a substitution rate at the lower end of

87 normal for single-stranded RNA viruses of between $1 \times 10^{-4}$ and $5 \times 10^{-4}$ substitutions per site

88    per year (16–18). This may be due to strong purifying selection (16), or due to peculiarities of

89    RABV. For example, the RABV genome is longer than average for RNA viruses, and genome

90    length and evolutionary rate are negatively correlated (19), although this relationship appears

91    to be weaker in single-stranded RNA viruses (20). A more unusual feature of RABV is that

92    infections can exhibit extended incubation periods within the host. The median generation

93    interval (the time between one individual becoming infected and then infecting another) is

94    estimated to be 17.3 days in domestic dogs (21), with other studies estimating mean serial

95    intervals of 26.3 days (22) and 45.0 days (23). Symptoms, infectivity, and death from rabies,

96    however, can occasionally occur years after the initial infection event (24). The length of the

97    incubation period is influenced by the route of exposure, with bites to the head and neck

98    leading to more rapid disease progression than bites to lower extremities (25). RABV can

99    remain in the muscle at the bite site for prolonged lengths of time before invading the host's

100   motor neurons and progressing through the nervous system, with limited, if any, infection of

101   other muscle fibres (26). While some replication in the muscle cells has been observed (27),

102   RABV replication at the inoculation site is not necessary for neural invasion (28). It is

103   currently unknown precisely how the RABV replication rate in the host muscle cells and

104   peripheral nervous system compares to the massive replication rate within the cells of the

105   central nervous system and brain. However, work suggests that RABV replication in muscle

106   cells may be reduced (29), and RABV replication in cultured rat sensory neurons may be 10-

107   to 100-fold lower than replication rates in rat and mouse CNS neurons (30). Rabies infections

108   that involve long incubation periods may, therefore, not lead to more accumulated mutations

109   than those with shorter incubation periods, as viral mutation is strongly influenced by the

110   replication process (31).

111   Changes in mutation rates through time due to long incubation periods may affect how we

112   analyse RABV sequence data and interpret these analyses. A relaxed molecular clock is

113    usually required to carry out phylogenetic analyses on rabies datasets, and it is not

114    uncommon for there to be difficulties in applying these analyses due to "insufficient temporal

115    signal"; usually referring to either no or a negative relationship between genetic divergence

116    and time, or this relationship having a lot of noise and a very low $R^2$ (32–36). RABV shows

117    variation in substitution rate between lineages (18,37,38) which may be driven in part by

118    differences in incubation periods. If the variable incubation period of rabies infections does

119    cause deviation from the molecular clock model (exceeding the variation captured by relaxed

120    or multirate clock models), this may negatively affect the accuracy of time-scaled

121    phylogenetic trees and emergence date predictions. Conversely, if mutation does continue at a

122    consistent rate during the incubation period, attention should be paid to extremely long

123    incubators which could drive the emergence of new variants, as seen recently in chronic

124    SARS-CoV-2 infections (39,40).

125    We hypothesised that reduced replication (and thus mutation) during the incubation period

126    could cause rabies evolution to be better represented by a per-generation model of mutation

127    than by the molecular clock model. We aim to clarify the nature of contemporary RABV

128    evolution using in silico methods, comparing the root-to-tip divergence of sequences

129    generated from synthetic outbreaks under per-unit time or per-generation mutation models,

130    and comparing these to RABV genomic data from Tanzania. We also aim to calculate a per-

131    generation substitution rate for RABV for future use as a parameter in transmission tree

132    reconstruction algorithms.

133

**Methods**

135    We investigate two contrasting mutational models for RABV – i.e., substitutions occurring on

136    a per-generation vs. per-unit-time basis – using a simulation approach. We first generated

137    synthetic RABV outbreaks using a branching process model (21) and then simulated these

138    two mutation processes over the resulting transmission trees. From the synthetic sequences

139    generated, we examined root-to-tip divergence and calculated variance explained ($R^2$) from

140    linear regressions, and compared these to the root-to-tip divergence of a set of RABV whole

141    genome sequences from Tanzania. Finally, we developed a method to estimate the per-

142    generation substitution rate for RABV and tested this on synthetic data before applying it to

143    the Tanzanian RABV dataset.

144

145    <u>Rabies outbreak simulation</u>

146    We simulated RABV mutation on branching-process simulations of rabies outbreaks.

147    Outbreaks were simulated 100 times over a spatially explicit representation of Mara Region

148    in northern Tanzania. In Serengeti District, where contact tracing data were available, the

149    model was initialised with the three cases that occurred in the mean generation interval ($g$=27

150    days, based on contact tracing data) prior to 2017 (simulations were run over a dog

151    population representing that in Mara region between 2017 and 2024). In the rest of Mara

152    region, where there were no data to guide initialisation, we seeded with ($0.01Dg$)/365 cases,

153    where $D$ is the initial dog population in that area. If $R_e$=1 (endemic transmission), this results

154    in roughly 1% of the population becoming rabid over a year; contact tracing data suggest that

155    incidence typically does not exceed that level (41). This led to a total of 273 initial cases in

156    the region. Each case was assigned a number of offspring cases drawn from a negative

157    binomial distribution (41) with mean ($R_0$)=1.05 and dispersion parameter=1.33. The $R_0$ value

158    was chosen to result in a median number of cases each month that was roughly constant over

159    time (over the 100 simulations), mimicking endemic disease. Movement of rabid dogs from

160    their home locations to and between transmission locations followed a random walk with step

161    lengths drawn from a Weibull distribution (shape=0.41, scale=0.13). We simulated occasional

162    long-distance transport of dogs to a random location prior to their first transmission in 2% of

163    cases (21). At each of a rabid dog's transmission locations, another dog was randomly

164    selected within the local 1km2 grid cell. If this dog was susceptible (i.e., not vaccinated or

165    already incubating infection from a prior transmission event), rabies was transmitted. A

166    generation interval was drawn for each new infection from a lognormal distribution

167    (meanlog=2.96, sdlog=0.82), describing the time delay before it also became rabid and made

168    its assigned transmissions. The step-length and generation-interval distributions were fitted

169    using contact tracing data from Serengeti District, Tanzania (21). Branching process

170    simulations were continued until 7 years had passed or rabies went extinct. Each synthetic

171    case was assigned an individual ID, and for every case (except initial seed cases) we recorded

172    the ID of the associated progenitor case. Dates of infection and transmission were recorded

173    for each case.

174    We isolated complete transmission trees descending from each of the 273 initial cases from

175    within one randomly selected synthetic outbreak. Transmission trees that contained over 100

176    cases (9 out of 273 trees in total, that ranged in size from 533 - 19,382 cases) were then used

177    to generate synthetic sequence data. Across these trees, we see a mean generation interval of

178    26.6 days, and 2.5 and 97.5 percentiles of 3.90 and 94.11 days (Supplementary Figure S1).

179    For each of the 9 trees the index case was assigned an initial 12kb genome sequence. Under

180    the per-unit time mutation model, we determined the expected number of mutations by

181    multiplying the substitution rate, the genome length and the length of the generation interval,

182    for each case along the resulting transmission tree (because we assume mutations are neutral,

183    the individual-level mutation rate is the same as the population-level substitution rate). The

184    realised number of mutations was then drawn from a Poisson distribution, with this mean. We

185    then randomly chose positions to change and new nucleotides to change them to. The

186 resulting synthetic sequence data is referred to as the "time-based sequence data". The

187 generation-based model of mutation works as above, with the exception that the expected

188 number of substitutions in a generation is constant and produces the synthetic "generation-

189 based sequence data".

190

191 <u>Divergence rate analysis</u>

192 To investigate patterns of temporal divergence under the mutation models described above,

193 we generated synthetic data with values of substitution rates ranging from 0.05 to 3

194 substitutions per genome per generation (or the per unit time substitution rate equivalent) and

195 4 population sampling regimes (from 1% of cases to 20%, informed by a previous study that

196 estimated that routine surveillance for rabies rarely confirms more than 10% of circulating

197 cases (42)). We calculated the genetic divergence as the number of nucleotide differences

198 from the index case to each sampled case. For each of the nine transmission trees, we then

199 compared genetic divergence with time under each scenario (substitution rate and sampling

200 regime combination), using linear regression through the origin.

201 In order to compare our synthetic patterns of divergence over time to real rabies data, a root-

202 to-tip divergence plot was also generated for a dataset of real RABV sequences (data from

203 (43); Figure 1A) using TempEst (v1.5.3 (13)), with the best-fit root located (Figure 1B).

204 These rabies cases occurred between 2001 and 2017 and were primarily from the Serengeti

205 district and Pemba Island, with the remaining sequences from elsewhere in Tanzania (Figure

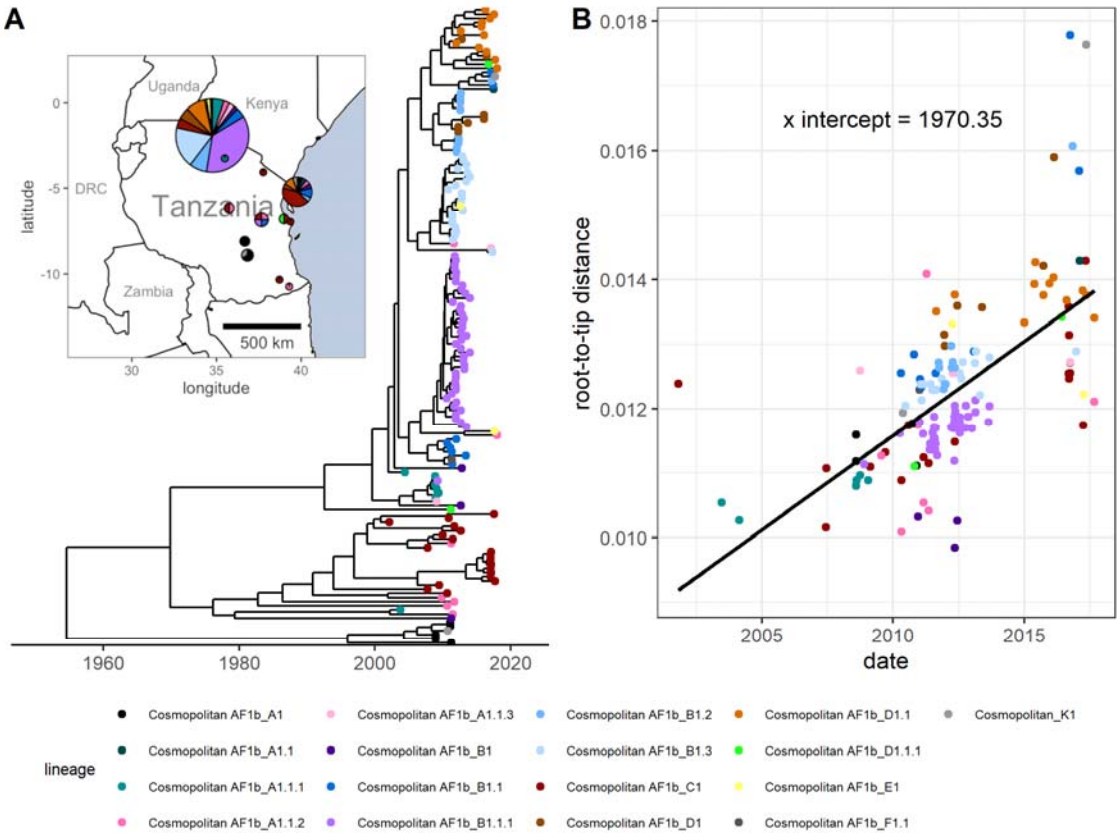206 1A inset). Sequence acquisition and tree building methods are detailed in (43).

207

208

**Fig 1. Phylogeny of real rabies virus whole genome sequences from Tanzania and root-to-tip divergence.** **(A)** The time-scaled tree (43) used to generate the root-to-tip divergence plot and to calculate the per-generation substitution rate. The inset map shows the approximate locations that the samples were collected from, and the lineages present in each location. Map point size represents the number of sequences in this dataset from district centroid locations. Base map data is from Natural Earth (naturalearthdata.com), via the *maps* R package. **(B)** The corresponding root-to-tip divergence plot. Point colours represent RABV lineage.

217

218     Calculating the per-generation substitution rate

219    We updated a method of calculating the per-generation substitution rate previously used in

220    eukaryotes (44) by using Bayesian posterior estimates of the clock rate and the generation

221    interval. We assessed this method's accuracy using the synthetic outbreak sequence data,

222    before applying it to the aforementioned set of RABV whole genome sequences.

223    To estimate the mean per-generation substitution rate, we analysed sequence data with

224    BEAST, and multiplied the posterior rate estimate for each MCMC sample (excluding the

225    burn-in period) by the generation-interval lengths sampled from the posterior of a simple

226    Bayesian analysis and then multiplied again by the genome length. The mean and 95%

227    credible interval of the estimate of the per-generation substitution rate for the RABV dataset

228    was calculated by taking the mean and the 2.5% and 97.5% percentiles of the resulting

229    multiplied posteriors.

230    To evaluate the accuracy of this method in estimating the mean per-generation substitution

231    rate, we also applied it to synthetic sequence data generated from outbreaks using the per-

232    generation mutation model as described above, under different substitution rates (11 values

233    ranging from 0.05 substitutions per generation to 3 substitutions per generation) and case

234    sampling rates (1%, 5%, 10% or 20% of cases sampled) across the 9 transmission trees that

235    contained at least 100 cases. Subsampled synthetic datasets containing more than 2000

236    sequences were not analysed as this number exceeds the total whole-genome RABV

237    sequences currently available on the RABV-GLUE database (45), and is unrealistic in the

238    context of examining individual rabies outbreaks. BEAST log files were generated from these

239    sequences using BEASTGen version 1.0.2 and BEAST version 1.10.4 (46). We chose to use a

240    JC substitution model with a strict clock, no site heterogeneity due to our per-generation

241    mutation model used in the simulations having equal probability of any site or base being

242    chosen and assumed constant population size. We used a tracelog frequency of 1000 and a

243    sufficiently long chain length for the effective sample size (ESS) of each parameter to exceed

244    200 when analysed using Tracer (47), and a 10% burn-in period. We applied the substitution

245    rate calculation method to these phylogenetic trees, and assessed the accuracy of the resulting

246    mean per-generation substitution rates by comparing them to the parameter values used to

247    generate the synthetic sequences, using the natural log of the ratio:

$$Deviation = \ln\left(\frac{M_e}{M_a}\right)$$

248    where $M_e$ is the mean estimated per-generation substitution rate and $M_a$ is the actual

249    substitution rate, where a deviation of zero means perfect accuracy.

250

251    The same method was applied to the dataset of 153 RABV sequences sampled from across

252    Tanzania (data from (43); Figure 1A). The mean per-generation substitution rate was

253    calculated, and distributions were fitted from the multiplied generation interval and clock rate

254    posteriors (generation interval posteriors based on values from (21) for the Tanzanian dataset,

255    extracted directly from the lognormal distribution used in simulations, and clock rate

256    posteriors taken from the BEAST log file of the time-scaled tree from Lushasi et al. (43)) and

257    genome length as described above. We compared different distributions (Gamma and

258    Lognormal) for estimating substitution rates and selected the best fitting distribution by AIC.

259    We also calculated the probabilities of between 0 and 10 SNP differences occurring across 1,

260    5 or 10 infection generations. For this calculation we simulated mutations arising at a Poisson

261    rate with lambda drawn from the fitted substitution rate distribution. The means and 95%

262    confidence intervals were calculated from the 10,000 simulations.

263

264    Data and code availability

265   All code is available at https://github.com/RowanDurrant/Rabies-Mutation. Analyses were

266   conducted using the R programming language (48). The beta regression curve and prediction

267   interval in Figure 2C was generated using the 'betareg' R package (49). RABV lineages were

268   assigned using MADDOG (45).

269

270   **Results**

271   <u>Root-to-tip divergence analysis</u>

272   At higher per-generation substitution rates (1 substitution per genome per generation and

273   above), distinct differences can be seen between root-to-tip divergence plots from the two

274   models of mutation (Figure 2A). The synthetic data generated from the per-generation

275   mutation model shows "stray" clusters or ridges of points both above and below the main

276   funnel of points, illustrated in the example in Figure 2A. Divergence plots from synthetic data

277   generated from the time-based model of mutation have less variance and do not exhibit this

278   pattern. At lower substitution rates (below 1 substitution per generation), no such pattern is

279   clearly distinguishable (Figure 2B). When the cases represented by the high-divergence

280   points from the per-generation model in Fig. 2A are visualised in a transmission tree, they are

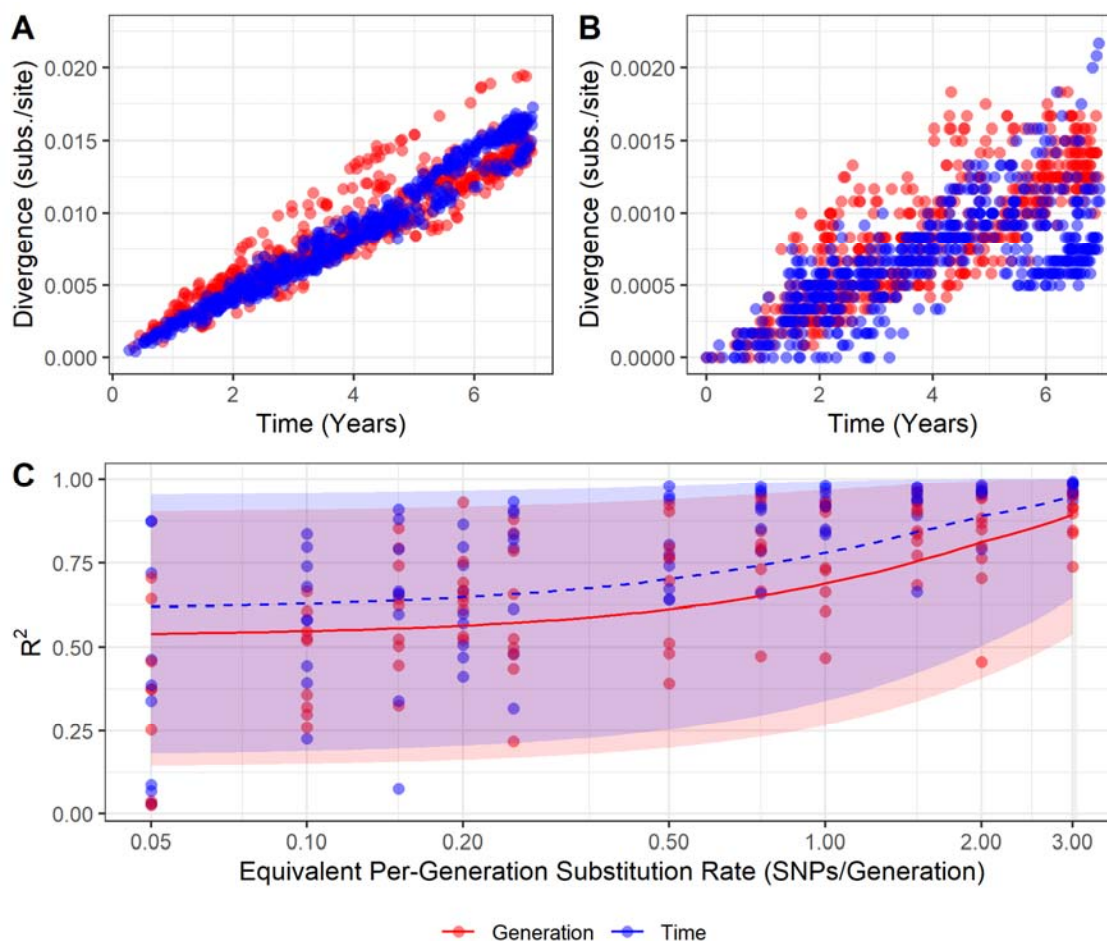281   mainly confined to a single chain (Supplementary Figure S1).

**Fig. 2: Temporal genetic divergence varies under two models of mutation. (A)** Root-to-tip divergence plots for synthetic sequences produced using time-based and generation-based mutation models, equivalent to 2 substitutions per genome per generation and **(B)** equivalent to 0.2 substitutions per genome per generation. Note that the y-axis scales differ by an order of magnitude between A and B. These data are from running mutation models over the same single transmission tree and have a case sampling rate of 5% (i.e., 621 cases sampled of 12,434 total). **(C)** The $R^2$ values obtained from regression through the origin of root-to-tip divergence of synthetic data from the time-based and generation-based models. Point colour indicates the mutation model used to generate the data. Lines represent beta regressions with logit links fit to data points, and shading represents the 95% prediction interval. The X axis is

293     log scaled. 5% of cases were sampled here; sampling rate had little effect on $R^2$

294     (Supplementary Figure S2).

295

296     Root-to-tip divergence plots derived from synthetic transmission trees using the time-based

297     mutation model had, on average, higher $R^2$ values than those from synthetic transmission

298     trees using the per-generation mutation model, although this is more difficult to distinguish

299     below a substitution rate of 0.5 substitutions per genome per generation (Figure 2C). As the

300     substitution rate increases, the $R^2$ values across both mutation models increase. The case

301     sampling rate appears to have little effect on $R^2$ (Supplementary Figure S2).

302     The root-to-tip divergence plot of the Tanzanian RABV dataset more closely resembles those

303     of lower substitution rate simulations, where it is difficult to determine any difference

304     between the models of mutation (Figure 1). While most lineages surround the regression line,

305     some (for example, Cosmopolitan AF1b_B1) group below the line, but without forming a

306     distinguishable "ridge".

307

308     <u>Substitution rate calculation</u>

309     The accuracy of our method used to calculate per-generation substitution rate remains similar

310     at all but the lowest values of substitution rate (Figure 3), with a tendency to underestimate

311     the substitution rate (meaning that the estimated substitution rate is below the substitution

312     rate parameter used to generate the synthetic data; mean natural log of the ratio of -0.18 and

313     root-mean-square of 0.54, where values of 0 indicate perfect estimates). Accuracy appears to

314     be more influenced by the number of sequences used in the BEAST analysis than by the case

315     sampling rate itself; the mean natural log of the ratio falls to -0.36 when fewer than 50

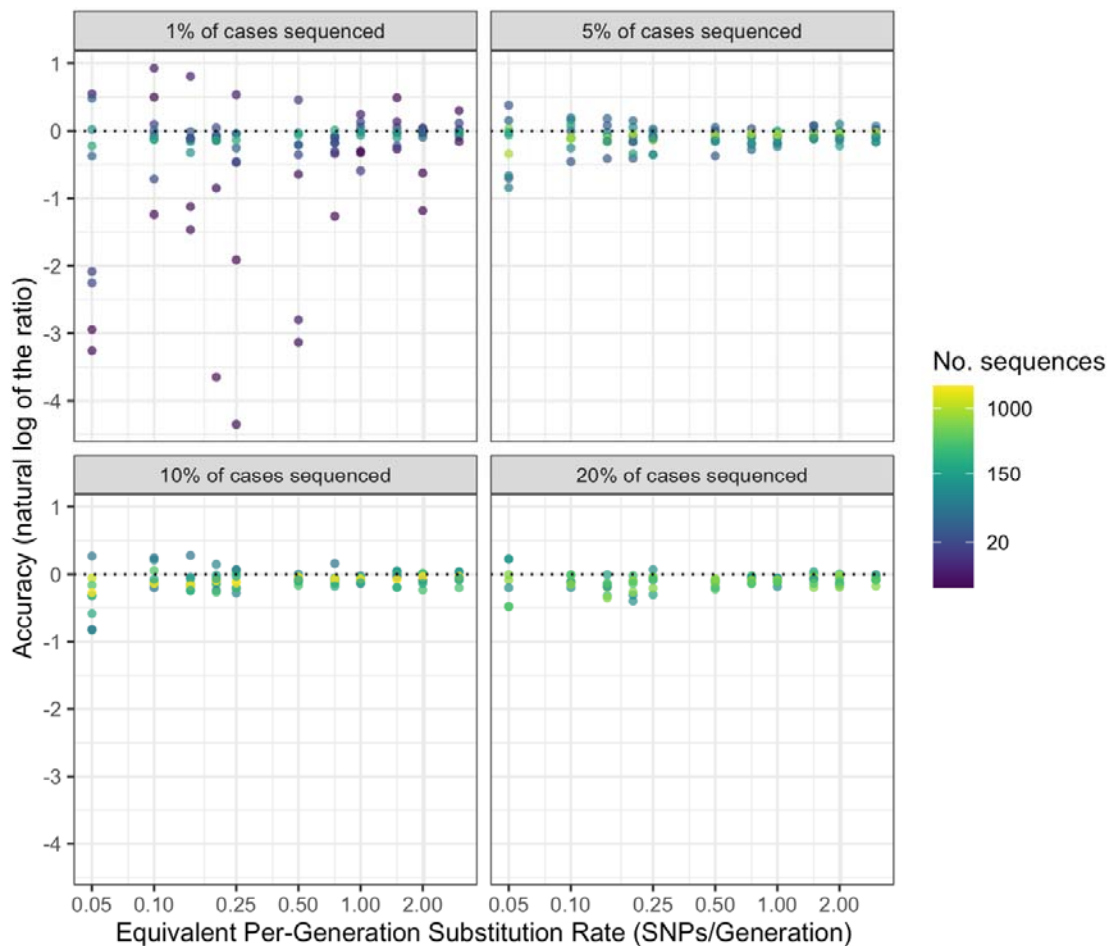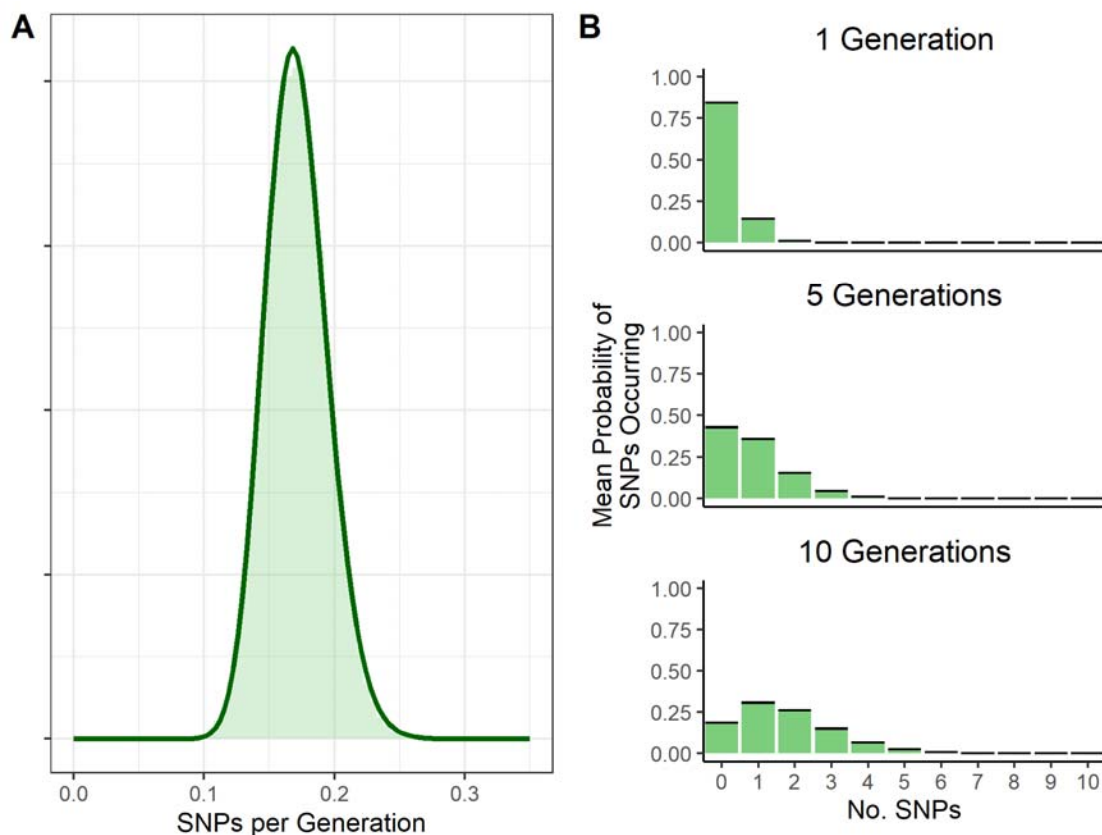316     sequences are used (root-square-mean of 0.95).

16

317



318

**Fig. 3: Accuracy of per-generation substitution rate predictions for different numbers of**

**sequences, substitution rates and sampling rates.** Facets indicate case sampling rate. The

dotted line represents perfect accuracy. X axis and colour scale are log transformed.

322

The Tanzanian RABV dataset from which we estimated the per-generation substitution rate

contains 153 sequences in total, and the accompanying time-scaled phylogenetic tree has a

root-to-tip height of approximately 65 years, although the sequences spanned just 16 years as

they were sampled from 2001 to 2017 (with 46.7% from years 2011-2012). These sequences

were largely complete; 98% of sequences were >95% complete (>11,327 kb in length). The

328    mean per generation substitution rate of RABV in this dataset was estimated to be 0.171

329    (95% credible interval: 0.127 - 0.219). The best fitting distribution by AIC to the output of

330    the multiplied Bayesian posteriors was a Gamma distribution with the parameters shape =

331    51.69 and rate = 301.8.

332

333    Using the calculated per generation per genome substitution rates, we calculated the

334    probability of different numbers of substitutions occurring over 1, 5 and 10 generations,

335    drawing the per-generation substitution rate ($\lambda$) from the above distribution (Figure 4). Over

336    many generations it is still quite likely for zero substitutions to occur; after 10 generations,

337    the probability of zero substitutions having occurred is 0.19.



338

339     **Fig. 4: Probability distributions of the mean per-generation substitution rate and**

340     **substitutions occurring over generations**. **(A)** estimated probability distribution of the per

341     genome per generation substitution rate from Tanzanian RABV sequences, with underlying

342     histogram of multiplied Bayesian posteriors of clock rate and generation interval. **(B)**

343     probability distribution of SNPs occurring over 1, 5 and 10 generations. The $\lambda$ value for a

344     Poisson rate of SNP occurrence is drawn from the SNPs per generation distribution fitted in

345     Figure 4A. Black bars represent the 95% confidence intervals (which are very tight).

346

347     **Discussion**

348     It can be difficult to get sufficient temporal signal for RABV sequence datasets, which we

349     hypothesised could be due in part to its variable incubation periods. We hypothesised that a

350     per-generation model of mutation may be more representative of RABV evolution than a

351     purely time-based model. We found that substantial differences in root-to-tip divergence

352     patterns between synthetic outbreaks using generation-based and time-based models of

353     mutation could be observed only at high underlying substitution rates. The substitution rate

354     for the Tanzanian RABV sequences examined (~0.17 substitutions per genome per

355     generation) was in the range where divergence patterns in the two models were extremely

356     similar. We can thus assume that the two models will give extremely similar results on the

357     relevant time scale. As we observed increasing divergence over time with reasonable $R^2$

358     values within this substitution rate range, it implies that variable incubation periods alone do

359     not fully account for the challenge in obtaining temporal signal. Therefore, other factors such

360     as insufficiently long sampling windows for the substitution rate are likely to be responsible

361     (50). This is an important consideration for analysing RABV sequences from new outbreaks,

362     or from endemic areas where sampling is opportunistic. As RABV has a substitution rate

363     lower than many other RNA viruses, longer sampling windows are required to achieve a

364     sufficient temporal signal.

365     The observation of little difference between root-to-tip divergence plots derived from the two

366     mutation models at substitution rates below 1 substitution per genome per generation is likely

367     because of averaging; multiple generations of infection are expected to have occurred per

368     substitution that arises on the viral genome. Over the many generations needed before

369     significant levels of viral genetic diversity are reached, the influence of any unusually long

370     incubation periods will be damped by the opposite influence of unusually short incubation

371     periods, eventually becoming indistinguishable from clock-like behaviour. On the other hand,

372     at higher substitution rates ridges form on the root-to-tip divergence plots under the per-

373     generation model of mutation but not under the per unit time model. While not affecting the

374     overall clock rate, these ridges reduce the overall $R^2$, and may be better analysed using a

375     separate local clock (51). The cases in these ridges almost all descend from a common

376     ancestor (Supplementary Figure S1), suggesting that a single unusually long or short

377     incubation period can affect which phylogenetic analyses we perform. Ridges caused by these

378     incubation periods can be distinguished from ridges caused by rate variation between lineages

379     as they will be parallel to the main cluster of points in the plot, whereas points belonging to

380     lineages with a different substitution rate will have a different slope. Studies examining the

381     number of substitutions occurring between successive sequenced cases, and whether this

382     increases when the secondary case's incubation period is unusually long, could clarify the

383     exact relationship between substitutions, generations, and time. More detailed data will be

384     required to investigate this further.

385     We calculated RABV's mean per-generation substitution rate to be approximately 0.17

386     substitutions per genome per transmission generation. This estimate is lower than those for

387     other RNA viruses, such as SARS (2 substitutions per genome per human passage (52)),

388    SARS-CoV-2 (0.52 substitutions per genome per 5.8-day generation interval (53)) and Ebola

389    virus (0.875 substitutions per genome per 14-day generation interval (54)). RNA viruses that

390    undergo periods of reduced replication or complete latency often show reduced substitution

391    rates, with one extreme example being HTLV-1/2 (55,56). However, we would not expect this

392    to affect the per-generation rate. The low per-generation substitution rate seen in rabies is

393    therefore likely due to mutation being constrained by other factors, such as strong purifying

394    selection (16), and likely contributes to the difficulties in obtaining sufficient temporal signal

395    for phylogenetic analyses. Previous studies suggest that for viruses in this substitution rate

396    range, sampling windows of up to 30 years may be required to overcome the phylodynamic

397    threshold (15); for comparison, SARS-CoV-2 achieved sufficient temporal signal within two

398    months of the start of the pandemic (50).

399    We can predict from the estimated per-generation substitution rate that identical sequences

400    are likely to have less than 5 intermediate generations between them (probability of fewer

401    than five generations occurring before a mutation occurs > 0.49 by repeated sampling of a

402    Poisson distribution with a lambda of 0.17), but still have a non-negligible probability of

403    being more distantly related. While the low substitution rate means that comparing the

404    number of SNPs between sequences alone may not be an effective method of determining

405    infector-infectee relationships, it could be used in conjunction with temporal and location

406    data to make more accurate predictions of transmission events by ruling out relationships

407    between more distantly related transmission chains co-circulating in the same area, as in (57).

408    Notably, our Poisson distribution of the number of substitutions occurring in one generation

409    is visually very similar to the genetic signature distribution reported in Cori *et al*. ((57), Fig

410    S1), despite different methods and RABV datasets being used in their calculations. It is likely,

411    however, that our estimate of the per-generation substitution rate is lower than the mean

412    number of SNPs expected between sequences from a primary and secondary case, due to the

413    time-based substitution rate being affected by purifying selection (58). Further analysis

414    comparing the estimated per-generation substitution rate to realised SNP distances between

415    primary-secondary case pairs could quantify this difference.

416    While the Jukes-Cantor model was the most appropriate to use on our synthetic data due to

417    the simplicity of the mutation models, phylogenetic analyses on real RABV genomes usually

418    use a more complex model, such as the GTR + G substitution model used to generate the

419    Tanzanian tree shown in this study (43). This, along with the simplicity of our mutation

420    model as well as sampling biases in the real dataset, may affect how comparable synthetic

421    root-to-tip divergence plots are to the real data.

422

423    While the molecular clock has proven critical for gaining insights into the history and

424    dynamics of disease outbreaks, the epidemiological characteristics of a virus should be

425    considered when choosing how to measure viral evolution. In this study, we determine that

426    the per-generation model is not likely to produce substantially different results from the

427    molecular clock model when analysing contemporary RABV evolution. We also estimate the

428    mean per-generation substitution rate of RABV for future use in transmission tree

429    reconstruction and efforts to estimate outbreak sizes and lineage emergence rates. Given that

430    many different lineages circulating simultaneously is seemingly a common occurrence in

431    areas with endemic rabies, it is important to investigate whether these lineages vary in

432    evolutionary rate and generation interval length, and ascertain the potential effects on

433    phylogenetic analyses.

434

435    **References**

436    1.    Gojobori T, Moriyama EN, Kimura M. Molecular clock of viral evolution, and the neutral theory.
437          Proc Natl Acad Sci. 1990 Dec;87(24):10015–8.

438    2.    Drummond AJ, Ho SYW, Phillips MJ, Rambaut A. Relaxed Phylogenetics and Dating with
439          Confidence. PLOS Biol. 2006 Mar 14;4(5):e88.

440    3.    Ho SYW, Duchêne S. Molecular-clock methods for estimating evolutionary rates and timescales.
441          Mol Ecol. 2014;23(24):5947–65.

442    4.    Drummond A, Oliver G. P, Rambaut A. Inference of Viral Evolutionary Rates from Molecular
443          Sequences. Adv Parasitol. 2003 Jan 1;54:331–58.

444    5.    Pybus OG, Rambaut A. Evolutionary analysis of the dynamics of viral infectious disease. Nat Rev
445          Genet. 2009 Aug;10(8):540–50.

446    6.    Wróbel B, Torres-Puente M, Jiménez N, Bracho MA, García-Robles I, Moya A, et al. Analysis of the
447          Overdispersed Clock in the Short-Term Evolution of Hepatitis C Virus: Using the E1/E2 Gene
448          Sequences to Infer Infection Dates in a Single Source Outbreak. Mol Biol Evol. 2006 Jun
449          1;23(6):1242–53.

450    7.    Grenfell BT, Pybus OG, Gog JR, Wood JLN, Daly JM, Mumford JA, et al. Unifying the
451          Epidemiological and Evolutionary Dynamics of Pathogens. Science. 2004 Jan 16;303(5656):327–
452          32.

453    8.    Gire SK, Goba A, Andersen KG, Sealfon RSG, Park DJ, Kanneh L, et al. Genomic surveillance
454          elucidates Ebola virus origin and transmission during the 2014 outbreak. Science. 2014 Sep
455          12;345(6202):1369–72.

456    9.    Kamath PL, Foster JT, Drees KP, Luikart G, Quance C, Anderson NJ, et al. Genomics reveals historic
457          and contemporary transmission dynamics of a bacterial disease among wildlife and livestock.
458          Nat Commun. 2016 May 11;7(1):11448.

459    10.   Drummond AJ, Pybus OG, Rambaut A, Forsberg R, Rodrigo AG. Measurably evolving populations.
460          Trends Ecol Evol. 2003 Sep 1;18(9):481–8.

461    11.   Korber B, Muldoon M, Theiler J, Gao F, Gupta R, Lapedes A, et al. Timing the Ancestor of the HIV-
462          1 Pandemic Strains. Science. 2000 Jun 9;288(5472):1789–96.

463    12.   Buonagurio DA, Nakada S, Parvin JD, Krystal M, Palese P, Fitch WM. Evolution of Human Influenza
464          A Viruses Over 50 Years: Rapid, Uniform Rate of Change in NS Gene. Science. 1986 May
465          23;232(4753):980–2.

466    13.   Rambaut A, Lam TT, Max Carvalho L, Pybus OG. Exploring the temporal structure of
467          heterochronous sequences using TempEst (formerly Path-O-Gen). Virus Evol. 2016
468          Jan;2(1):vew007.

469    14.   Duchene S, Lemey P, Stadler T, Ho SYW, Duchene DA, Dhanasekaran V, et al. Bayesian Evaluation
470          of Temporal Signal in Measurably Evolving Populations. Mol Biol Evol. 2020 Nov 1;37(11):3363–
471          79.

472    15.   Duchêne S, Duchêne D, Holmes EC, Ho SYW. The Performance of the Date-Randomization Test in
473          Phylogenetic Analyses of Time-Structured Virus Data. Mol Biol Evol. 2015 Jul 1;32(7):1895–906.

474     16.  Holmes EC, Woelk CH, Kassis R, Bourhy H. Genetic Constraints and the Adaptive Evolution of
475           Rabies Virus in Nature. Virology. 2002 Jan 20;292(2):247–57.

476     17.  Biek R, Pybus OG, Lloyd-Smith JO, Didelot X. Measurably evolving pathogens in the genomic era.
477           Trends Ecol Evol. 2015 Jun 1;30(6):306–13.

478     18.  Layan M, Dellicour S, Baele G, Cauchemez S, Bourhy H. Mathematical modelling and
479           phylodynamics for the study of dog rabies dynamics and control: A scoping review. PLoS Negl
480           Trop Dis. 2021 May 27;15(5):e0009449.

481     19.  Duchêne S, Holmes EC. Estimating evolutionary rates in giant viruses using ancient genomes.
482           Virus Evol. 2018 Feb 27;4(1):vey006.

483     20.  Sanjuán R. From Molecular Genetics to Phylodynamics: Evolutionary Relevance of Mutation
484           Rates Across Viruses. PLoS Pathog. 2012 May 3;8(5):e1002685.

485     21.  Mancy R, Rajeev M, Lugelo A, Brunker K, Cleaveland S, Ferguson EA, et al. Rabies shows how
486           scale of transmission can enable acute infections to persist at low prevalence. Science. 2022 Apr
487           29;376(6592):512–6.

488     22.  Hayes S, Lushasi K, Sambo M, Changalucha J, Ferguson EA, Sikana L, et al. Understanding the
489           incidence and timing of rabies cases in domestic animals and wildlife in south-east Tanzania in
490           the presence of widespread domestic dog vaccination campaigns. Vet Res. 2022 Dec
491           12;53(1):106.

492     23.  Kurosawa A, Tojinbara K, Kadowaki H, Hampson K, Yamada A, Makita K. The rise and fall of rabies
493           in Japan: A quantitative history of rabies epidemics in Osaka Prefecture, 1914–1933. PLoS Negl
494           Trop Dis. 2017 Mar 23;11(3):e0005435.

495     24.  Boland TA, McGuone D, Jindal J, Rocha M, Cumming M, Rupprecht CE, et al. Phylogenetic and
496           epidemiologic evidence of multiyear incubation in human rabies. Ann Neurol. 2014;75(1):155–
497           60.

498     25.  Dimaano EM, Scholand SJ, Alera MTP, Belandres DB. Clinical and epidemiological features of
499           human rabies cases in the Philippines: a review from 1987 to 2006. Int J Infect Dis. 2011 Jul
500           1;15(7):e495–9.

501     26.  Charlton KM, Nadin-Davis S, Casey GA, Wandeler AI. The long incubation period in rabies:
502           delayed progression of infection in muscle at the site of exposure. Acta Neuropathol (Berl). 1997
503           Jun 1;94(1):73–7.

504     27.  Yamaoka S, Ito N, Ohka S, Kaneda S, Nakamura H, Agari T, et al. Involvement of the Rabies Virus
505           Phosphoprotein Gene in Neuroinvasiveness. J Virol. 2013 Nov 15;87(22):12327–38.

506     28.  Shankar V, Dietzschold B, Koprowski H. Direct entry of rabies virus into the central nervous
507           system without prior local replication. J Virol. 1991 May;65(5):2736–8.

508     29.  Schnell MJ, McGettigan JP, Wirblich C, Papaneri A. The cell biology of rabies virus: using stealth
509           to reach the brain. Nat Rev Microbiol. 2010 Jan;8(1):51–61.

510     30.  Lycke E, Tsiang H. Rabies virus infection of cultured rat sensory neurons. J Virol. 1987
511           Sep;61(9):2733–41.

512   31.  Belshaw R, Gardner A, Rambaut A, Pybus OG. Pacing a small cage: mutation and RNA viruses.
513        Trends Ecol Evol. 2008 Apr 1;23(4):188–93.

514   32.  Fusaro A, Monne I, Salomoni A, Angot A, Trolese M, Ferrè N, et al. The introduction of fox rabies
515        into Italy (2008–2011) was due to two viral genetic groups with distinct phylogeographic
516        patterns. Infect Genet Evol. 2013 Jul 1;17:202–9.

517   33.  Wang L, Wu X, Bao J, Song C, Du J. Phylodynamic and transmission pattern of rabies virus in
518        China and its neighboring countries. Arch Virol. 2019 Aug 1;164(8):2119–29.

519   34.  Zhang Y, Vrancken B, Feng Y, Dellicour S, Yang Q, Yang W, et al. Cross-border spread, lineage
520        displacement and evolutionary rate estimation of rabies virus in Yunnan Province, China. Virol J.
521        2017 Jun 3;14(1):102.

522   35.  Faye M, Faye O, Paola ND, Ndione MHD, Diagne MM, Diagne CT, et al. Rabies surveillance in
523        Senegal 2001 to 2015 uncovers first infection of a honey-badger. Transbound Emerg Dis.
524        2022;69(5):e1350–64.

525   36.  Caraballo DA, Lema C, Novaro L, Gury-Dohmen F, Russo S, Beltrán FJ, et al. A Novel Terrestrial
526        Rabies Virus Lineage Occurring in South America: Origin, Diversification, and Evidence of Contact
527        between Wild and Domestic Cycles. Viruses. 2021 Dec;13(12):2484.

528   37.  Troupin C, Dacheux L, Tanguy M, Sabeta C, Blanc H, Bouchier C, et al. Large-Scale Phylogenomic
529        Analysis Reveals the Complex Evolutionary History of Rabies Virus in Multiple Carnivore Hosts.
530        PLOS Pathog. 2016 Dec 15;12(12):e1006041.

531   38.  Streicker DG, Lemey P, Velasco-Villa A, Rupprecht CE. Rates of Viral Evolution Are Linked to Host
532        Geography in Bat Rabies. PLOS Pathog. 2012 May 17;8(5):e1002720.

533   39.  Kemp SA, Collier DA, Datir RP, Ferreira IATM, Gayed S, Jahun A, et al. SARS-CoV-2 evolution
534        during treatment of chronic infection. Nature. 2021 Apr;592(7853):277–82.

535   40.  Choi B, Choudhary MC, Regan J, Sparks JA, Padera RF, Qiu X, et al. Persistence and Evolution of
536        SARS-CoV-2 in an Immunocompromised Host. N Engl J Med. 2020 Dec 3;383(23):2291–3.

537   41.  Hampson K, Dushoff J, Cleaveland S, Haydon DT, Kaare M, Packer C, et al. Transmission Dynamics
538        and Prospects for the Elimination of Canine Rabies. PLOS Biol. 2009 Mar 10;7(3):e1000053.

539   42.  Townsend SE, Lembo T, Cleaveland S, Meslin FX, Miranda ME, Putra AAG, et al. Surveillance
540        guidelines for disease elimination: A case study of canine rabies. Comp Immunol Microbiol Infect
541        Dis. 2013 May;36(3):249–61.

542   43.  Lushasi K, Brunker K, Rajeev M, Ferguson EA, Jaswant G, Baker LL, et al. Integrating contact
543        tracing and whole-genome sequencing to track the elimination of dog-mediated rabies: an
544        observational and genomic study. Flegg J, editor. eLife. 2023 May 25;12:e85262.

545   44.  Slatkin M, Hudson RR. Pairwise comparisons of mitochondrial DNA sequences in stable and
546        exponentially growing populations. Genetics. 1991 Oct 1;129(2):555–62.

547   45.  Campbell K, Gifford RJ, Singer J, Hill V, O'Toole A, Rambaut A, et al. Making genomic surveillance
548        deliver: A lineage classification and nomenclature system to inform rabies elimination. PLOS
549        Pathog. 2022 May 2;18(5):e1010023.

550   46.  Suchard MA, Lemey P, Baele G, Ayres DL, Drummond AJ, Rambaut A. Bayesian phylogenetic and
551         phylodynamic data integration using BEAST 1.10. Virus Evol. 2018 Jan 1;4(1):vey016.

552   47.  Rambaut A, Drummond AJ, Xie D, Baele G, Suchard MA. Posterior Summarization in Bayesian
553         Phylogenetics Using Tracer 1.7. Syst Biol. 2018 Sep 1;67(5):901–4.

554   48.  R Core Team. R: A Language and Environment for Statistical Computing [Internet]. Vienna,
555         Austria; 2020. Available from: https://www.R-project.org/

556   49.  Cribari-Neto F, Zeileis A. Beta Regression in R. J Stat Softw. 2010 Apr 5;34(1):1–24.

557   50.  Duchene S, Featherstone L, Haritopoulou-Sinanidou M, Rambaut A, Lemey P, Baele G. Temporal
558         signal and the phylodynamic threshold of SARS-CoV-2. Virus Evol. 2020 Jul 1;6(2):veaa061.

559   51.  Featherstone LA, Rambaut A, Duchene S, Wirth W. Clockor2: Inferring global and local strict
560         molecular clocks using root-to-tip regression [Internet]. bioRxiv; 2023 [cited 2023 Aug 17]. p.
561         2023.07.13.548947. Available from:
562         https://www.biorxiv.org/content/10.1101/2023.07.13.548947v1

563   52.  Vega VB, Ruan Y, Liu J, Lee WH, Wei CL, Se-Thoe SY, et al. Mutational dynamics of the SARS
564         coronavirus in cell culture and human populations isolated in 2003. BMC Infect Dis. 2004 Sep
565         6;4:32.

566   53.  Braun K, Moreno G, Wagner C, Accola MA, Rehrauer WM, Baker D, et al. Limited within-host
567         diversity and tight transmission bottlenecks limit SARS-CoV-2 evolution in acutely infected
568         individuals [Internet]. bioRxiv; 2021 [cited 2023 Feb 9]. p. 2021.04.30.440988. Available from:
569         https://www.biorxiv.org/content/10.1101/2021.04.30.440988v1

570   54.  Kinganda-Lusamaki E, Black A, Mukadi D, Hadfield J, Mbala-Kingebeni P, Pratt CB, et al.
571         Operationalizing genomic epidemiology during the Nord-Kivu Ebola outbreak, Democratic
572         Republic of the Congo [Internet]. medRxiv; 2020 [cited 2023 Feb 9]. p. 2020.06.08.20125567.
573         Available from: https://www.medrxiv.org/content/10.1101/2020.06.08.20125567v1

574   55.  Holmes EC. Molecular Clocks and the Puzzle of RNA Virus Origins. J Virol. 2003 Apr;77(7):3893–7.

575   56.  Van Dooren S, Salemi M, Vandamme AM. Dating the Origin of the African Human T-Cell
576         Lymphotropic Virus Type-I (HTLV-I) Subtypes. Mol Biol Evol. 2001 Apr 1;18(4):661–71.

577   57.  Cori A, Nouvellet P, Garske T, Bourhy H, Nakouné E, Jombart T. A graph-based evidence synthesis
578         approach to detecting outbreak clusters: An application to dog rabies. PLOS Comput Biol. 2018
579         Dec 17;14(12):e1006554.

580   58.  Duchêne S, Holmes EC, Ho SYW. Analyses of evolutionary dynamics in viruses are hindered by a
581         time-dependent bias in rate estimates. Proc R Soc B Biol Sci. 2014 Jul 7;281(1786):20140732.

582

583   **Supporting information captions**

584 **Supp. Fig. S1: histogram of generation intervals from the simulated outbreaks.** Vertical

585 dashed lines represent the median (blue) and mean (red) generation interval.

586 **Supp. Fig. S2: points in the offshoot ridge predominantly occur in one transmission tree.**

587 **(A)** root-to-tip divergence plot (2 SNPs/genome/generation, 5% of cases sequenced) with

588 offshoot ridge points highlighted in red. Offshoot ridge points are defined in this plot as

589 having a divergence rate above $8\times10^{-6}$ substitutions/day and occurring after day 750. **(B)**

590 transmission tree of the simulated outbreak with offshoot ridge cases highlighted in red.

591 Graph edge length is not proportional to time or divergence.

592 **Supp. Fig. S3: Sampling rate does not impact the $R^2$ of root-to-tip divergence plots from**

593 **synthetic data.** Plot is faceted by the proportion of the total number of cases in the outbreak

594 sequenced, point colour represents mutation model.

595