

1 **Current sampling and sequencing biases of Lassa mammarenavirus limit**
2 **inference from phylogeography and molecular epidemiology in Lassa Fever**
3 **endemic regions.**

4
5 **Authors**

6 Liã Bárbara Arruda^{1b#}, Hayley Beth Free^{2a§}, David Simons^{2§}, Rashid Ansumana³, Linzy Elton¹,
7 Najmul Haider^{2c}, Isobella Honeyborne¹, Danny Asogun⁴, Timothy D McHugh¹, Francine
8 Ntoumi^{5,6}, Alimuddin Zumla^{1,7}, Richard Kock²

9
10 **Affiliations**

11 ¹ Centre for Clinical Microbiology, Division of Infection and Immunity, University College
12 London, London, UK

13 ² The Royal Veterinary College, University of London, Hatfield, UK.

14 ³ School of Community Health Sciences, Njala University, Bo, Sierra Leone

15 ⁴ Ekpoma and Irrua Specialist Teaching Hospital, Ambrose Alli University, Irrua, Nigeria.

16 ⁵ Fondation Congolaise pour la Recherche Médicale (FCRM), Brazzaville, Republic of Congo

17 ⁶ Institute for Tropical Medicine, University of Tübingen, Germany

18 ⁷ NIHR Biomedical Research Centre, UCL Hospitals NHS Foundation Trust, London, UK

19 ^a Current affiliation Oxford Brookes University, Oxford, UK

20 ^b Current affiliation Wellcome Connecting Science, Hinxton, UK

21 ^c Current affiliation School of Life Sciences, Faculty of Natural Sciences, Keele University,
22 Staffordshire, United Kingdom

23 [§] Both authors contributed equality to this work

24 [#] Corresponding author

27 **Abstract**

28 Lassa fever (LF) is a potentially lethal viral haemorrhagic infection of humans caused by *Lassa*
29 *mammarenavirus* (LASV). It is an important endemic zoonotic disease in West Africa with
30 growing evidence for increasing frequency and sizes of outbreaks. Phylogeographic and
31 molecular epidemiology methods have projected expansion of the Lassa fever endemic zone
32 in the context of future global change. The Natal multimammate mouse (*Mastomys natalensis*)
33 is the predominant LASV reservoir, with few studies investigating the role of other animal
34 species. To explore host sequencing biases, all LASV nucleotide sequences and associated
35 metadata available on GenBank (n = 2,298) were retrieved. Most data originated from Nigeria
36 (54%), Guinea (20%) and Sierra Leone (14%). Data from non-human hosts (n = 703) were
37 limited and only 69 sequences encompassed complete genes. We found a strong positive
38 correlation between the number of confirmed human cases and sequences at the country level
39 ($r = 0.93$ (95% Confidence Interval = 0.71 - 0.98), $p < 0.001$) but no correlation exists between
40 confirmed cases and the number of available rodent sequences ($r = -0.019$ (95% C.I. -0.71 -
41 0.69), $p = 0.96$). Spatial modelling of sequencing effort highlighted current biases in locations
42 of available sequences, with increased effort observed in Southern Guinea and Southern
43 Nigeria. Phylogenetic analyses showed geographic clustering of LASV lineages, suggestive
44 of isolated events of human-to-rodent transmission and the emergence of currently circulating
45 strains of LASV from the year 1498 in Nigeria. Overall, the current study highlights significant
46 geographic limitations in LASV surveillance, particularly, in non-human hosts. Further
47 investigation of the non-human reservoir of LASV, alongside expanded surveillance, are
48 required for precise characterisation of the emergence and dispersal of LASV. Accurate
49 surveillance of LASV circulation in non-human hosts is vital to guide early detection and
50 initiation of public health interventions for future Lassa fever outbreaks.

51

52 **Key-words**

53 Lassa mammarenavirus; Lassa Fever; Phylogeography; Metadata; Zoonoses; Surveillance

54

55

56 **1 Introduction**

57

58 Lassa fever (LF) is a lethal zoonotic viral haemorrhagic disease of humans, caused by *Lassa*
59 *mammarenavirus* (LASV). It causes an estimated 900,000 annual human infections and
60 several thousand deaths in West Africa annually (1,2). The WHO assigns LASV endemicity to
61 eight West African countries: Benin, Ghana, Guinea, Liberia, Mali, Sierra Leone, Togo and
62 Nigeria (S1 Fig) (3). LASV is a bisegmented ssRNA- virus of the family Arenaviridae (4,5).
63 Based on the genomic analysis of the large (L) and small segments (S) LASV has been
64 classified into seven lineages which demonstrate spatial segregation across the endemic
65 range (6). The high nucleotide variability (25-32%) of these lineages introduces complexity
66 into assays to detect LASV infection.

67

68 Epidemiological data on LF is limited and constrained by current testing and reporting in the
69 endemic region, making accurate estimates of its true burden challenging (7). Many individuals
70 infected with LASV do not seek healthcare with up to 80% of infections assumed
71 asymptomatic or presenting as mild illness (8). Estimates based on longitudinal serological
72 surveys in Sierra Leone in the early 1980's indicated that 100,000 to 300,000 infections of LF
73 occurred annually in West Africa, with more recent estimates being up to 900,000 infections
74 (2,8). Identification of symptomatic cases is further confounded by overlapping symptoms with
75 other diseases (e.g., malaria) and lack of available diagnostic methods (1,9–11). Access to
76 diagnostic tests varies spatially, increased availability at centers of excellence in LF treatment
77 and research such as the Irrua Specialist Teaching Hospital, Nigeria and Kenema General
78 Hospital, Sierra Leone results in a spatial bias of reported cases from these locations.
79 Phylogenetic analysis and molecular dating of sequence clinical and research samples
80 suggest a westward route of dispersal of LASV lineages, from the most recent common
81 ancestor in Nigeria. (12–18). These estimates have been used to project the potential for
82 Lassa Fever to extend beyond the current endemic zone (19).

83

84 The Natal multimammate mouse (*Mastomys natalensis*) is the primary reservoir of LASV,
85 however, 11 other rodent species have been found to be acutely infected or have seropositivity
86 to LASV including; *Mastomys erythroleucus*, *Hylomyscus pamfi*, *Mus baoulei* and *Rattus*
87 *rattus* (15,20–24). Humans become infected with LASV upon contact with or inhalation of
88 excretions from the rodent species (12,25). Although human-to-human transmission has been
89 reported – typically associated with nosocomial outbreaks – these are rare events when
90 compared with spillover from rodent hosts (26). We performed a study of LASV nucleotide
91 sequences available from the National Centre for Biotechnology Information (NCBI) GenBank,
92 using associated metadata to spatially model sequencing effort, adjusted for the number of

93 suspected and confirmed human LF cases to determine potential biases in locations of
94 available sequences or significant geographic limitations in LASV surveillance, particularly, in
95 non-human hosts.

96

97

98 **2 Methods**

99

100 **2.1 Data Collection and Processing**

101

102 LASV nucleotide and protein sequences were obtained from the NCBI GenBank (27). The
103 search query run on 24 Sep 2021 was for “Lassa mammarenavirus” in the organism field of
104 the NCBI nucleotide dataset. Data were obtained using the NCBI Entrez API with analysis
105 conducted using the “genbankr” package within the R statistical programming language (27–
106 29). Associated citations were manually retrieved to identify missing metadata for sequences
107 including hosts and geographic location of samples. Sequences with large portions (10%
108 missing compared to reference sequences, NC_004296.1 and NC_004297.1 for S and L
109 segments respectively) of missing nucleotide data on the L- or S-segment or lacking
110 associated metadata (collection year, host species, country, and geographical region of
111 sampling) were excluded from phylogenetic analysis. Nucleotide sequences were aligned
112 using the ‘map to reference’ tool on Geneious Prime 20201.2. Alignment, visual inspection
113 and manual editing were performed, and entries that contained >100 continuous ambiguous
114 nucleotide calls were excluded (S1 Data).

115

116 **2.2 Sequencing Bias**

117

118 First, we compared the number of cases reported from countries between 2008-2023 with the
119 number of samples contained in GenBank to summarise the correlation between reported
120 human cases and availability of sequences. We then compared the proportion of human to
121 non-human derived sequences within countries.

122

123 To understand the bias of sequenced samples at a sub-national level the origin of a sequenced
124 sample was geocoded using the Google Geocoding API using the “ggmap” package (30).
125 Sequence locations were associated with level-1 administrative regions and data were
126 separated into human and rodent sources of samples to visualise the spatial heterogeneity of
127 sampling. To measure sampling effort bias, the number of samples obtained within a level-1
128 administrative region was associated with the centroid of the region. The number of confirmed
129 LF clinical cases reported from these regions in the previous 15 years was obtained (S2 Data).

130 The number of cases within a region was divided by the human population count to produce
131 the number of confirmed cases per 100,000 individuals. The number of sequences was used
132 as the response variable in a spatial Generalised Additive Model, with geographic coordinates
133 and cases per 100,000 individuals used as covariates. This model was constructed using the
134 “mgcv” package (31).

135

136 2.3 Phylogenetic Analysis

137

138 Phylogenetic analysis was undertaken through Bayesian Markov Chain Monte Carlo (MCMC)
139 method using BEAST.v1.10.4 (32). In BEAUTi, the parameters were a substitution model as
140 a generalised time reversible plus gamma site heterogeneity, with codon partition positions 1,
141 2, 3. A strict clock and a coalescent tree prior with a constant size population was used. Each
142 analysis consisted of 20 million MCMC steps and trees were sampled every 20,000
143 generations. Sample collection dates from the metadata were used as tip dates to fit to a
144 molecular clock, and country of sample collection was incorporated as a discrete state (16,33).
145 To assess the log files of the output TRACER.v.1.7.1 was used. Maximum-clade credibility
146 trees were generated through TreeAnnotator v1.8.4 and visualised in FigTree.v1.4.4 (34).

147

148 3 Results

149

150 3.1 Compiled Dataset

151

152 The initial dataset comprised 2,298 records (from samples obtained 1969-2019), including
153 nucleotide sequences and associated metadata. Incomplete gene sequences and sequences
154 lacking metadata information (n = 1,045) were removed from phylogenetic analyses.
155 Therefore, 680 sequences of complete S segment and 573 sequences of partial L segment (L
156 protein only) were used. Accession numbers of included and excluded sequences are
157 available in S1 Data.

158

159 3.2 Descriptive Analysis

160

161 Year of collection was available for 2,108 records, with the oldest sequence dating from 1969
162 and latest from 2019. Among these records, most sequences (n = 1,936, 92%) have been
163 obtained since 2008. Human-derived LASV sequences comprised most of the available
164 records (67%), other host species include *Mastomys natalensis* (29%) and *Mastomys spp.*
165 (3%), while *Mastomys erythroleucus* (n = 18), *Mus baoulei* (n = 9) and *Hylomyscus pamfi* (n =
166 10) represent < 1% each. The species sampled was not documented in 107 records. Country

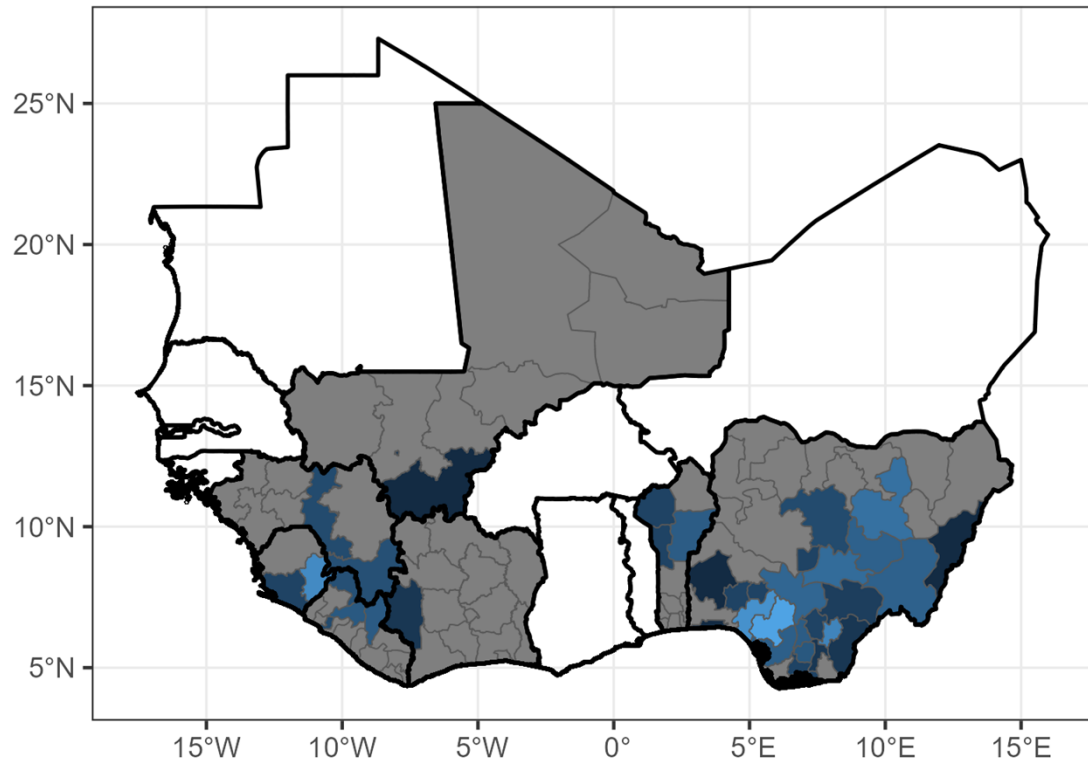
167 of collection was available for 2,238 records. Most sequences were produced from samples
168 collected in Nigeria (54%), followed by Guinea (20%), Sierra Leone (14%), Liberia (4%) and
169 Cote d'Ivoire (3%) with the remainder obtained from, Benin, Ghana, Mali and Togo (Fig 1).

170

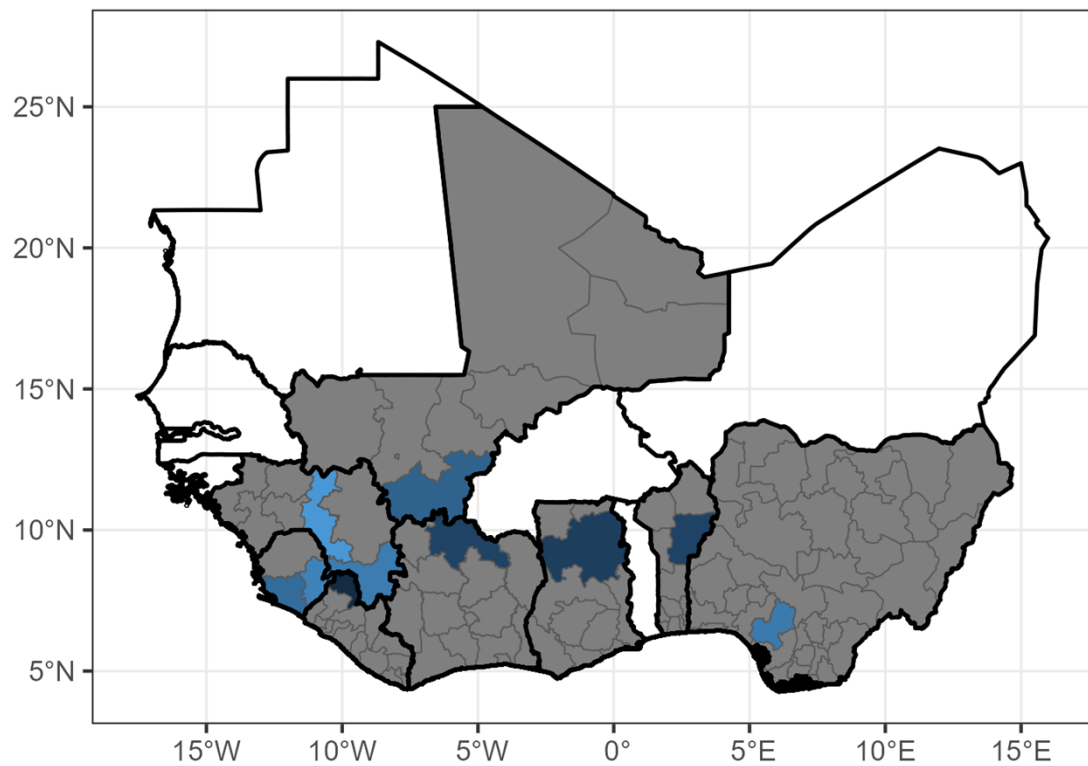
171 Sequences for human derived samples with regional location data (n = 1328, 63%) were
172 clustered in Edo State, Nigeria (n = 519, 39%), Ondo State, Nigeria (n = 220, 17%) and
173 Eastern Province, Sierra Leone (n = 159, 12%) with 430 samples from the remaining endemic
174 regions. Sequences from rodent samples with regional location data (n = 527, 25%) were most
175 commonly obtained from Faranah, Guinea (n = 210, 39%) and Eastern Province, Sierra Leone
176 (n = 107, 20%) with 210 samples from the regions.

177

Human



Rodent



Number of sequences (log₁₀)

0 1 2 3

179 Figure 1 – The number of sequences, shown on a \log_{10} scale, retrieved from NCBI GenBank
180 with associated regional sampling location and host for human samples (top, $n = 1,328$) and
181 rodent samples (bottom, $n = 527$). Grey regions represent level-1 administrative areas with no
182 sequences within countries that have at least one available sequence. White countries are
183 West African countries with no available LASV sequences. See S1 Fig for country names.
184 Shapefiles for basemap layer obtained from GADM 4.0.2 (35)

185
186

187 3.3 Sequencing bias

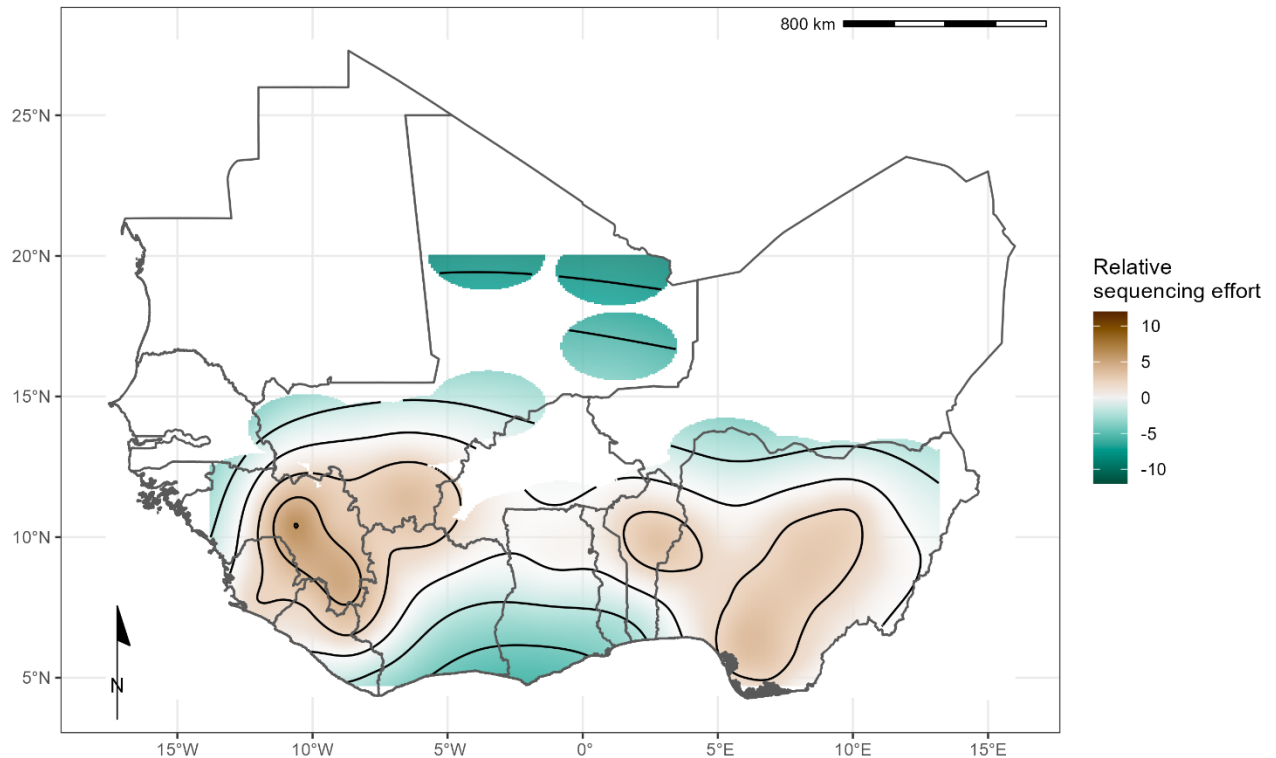
188

189 We observed a strong positive correlation between the number of confirmed human cases
190 between 2008-2023 and the number of GenBank deposited sequences at country level
191 ($r(\text{degrees of freedom} = 7) = 0.93$ (95% Confidence Interval = 0.71-0.98), $p < 0.001$). When
192 analysed by species source no correlation was observed with the number of confirmed cases
193 and the number of available rodent sequences was observed ($r(6) = -0.019$ (95% C.I. -0.71-
194 0.69), $p = 0.96$).

195

196 When combining both human and rodent-derived samples at the regional level to explore
197 spatial sampling biases, we found that sequencing effort is greatest in Southwest Nigeria,
198 centred over Edo State and the Faranah and Nzérékoré regions of Guinea, Eastern Province
199 of Sierra Leone and Nimba district of Liberia (Fig 2). There was a positive, non-linear
200 association between the rate of confirmed human cases with the number of available rodent
201 and human derived LASV sequences at regional level (deviance explained = 14%, estimated
202 degrees of freedom = 2.29, $p < 0.001$).

203



204

205 Figure 2 – Modelled relative sequencing effort derived from both human and rodent samples.
206 Greatest sequencing effort coincides with areas where sampling in humans (Edo, Nigeria and
207 Kenema, Sierra Leone) and rodents (Faranah, Guinea) have historically been focused.
208 Shapefiles for basemap layer obtained from GADM 4.0.2 (35)

209

210 3.4 Phylogenetic Analysis

211

212 Sequences for each segment of LASV showed clustering according to previously documented
213 lineages I-VII alongside geographical clustering with lineages I-III and VI present in Nigeria,
214 IV in Liberia, Guinea and Sierra Leone, V in Mali and VII in Togo (S2 Fig). In this analysis only
215 L segment sequences of lineage V from Cote d'Ivoire were included due to quality control
216 exclusion criteria. The phylogeny of the L segment indicates an older emergence of LASV in
217 the human population, with the most recent common ancestor (MRCA) predicted in the year
218 828 in Nigeria, inference based on the S segment indicates the emergence in the year 1350
219 (Table 1).

220

221 Table 1 - The most recent common ancestor (MRCA) stratified by host and country of
 222 collection of *Lassa marmarenavirus* (LASV) S and L segments. Samples were collected
 223 between 1969-2018.

Host species	Country	S segment MRCA	L segment MRCA
<i>Homo sapiens</i> (n=1181)	Benin	1995	1989
	Guinea	1895	1871
	Liberia	1895	1627
	Nigeria	1681	1498
	Sierra Leone	1901	1874
	Togo	2016	2014
<i>Hylomyscus pamfi</i> (n=2)	Nigeria	1681	1498
<i>Mastomys erythroleucus</i> (n=18)	Guinea	1975	2010
	Nigeria	2008	2006
<i>Mastomys natalensis</i> (n=36)	Guinea	1938	1997
	Mali	1951	2007
	Sierra Leone	1909	1979

224
 225 There was a lack of sequence information from lineage I and VI, however, phylogeny suggests
 226 these lineages are basal to others in Nigeria (S2 Fig). Lineage VII in Togo is most closely
 227 related to Nigerian isolates and potentially diverged between 500-900 years ago. The
 228 divergence of lineage III and IV is predicted to have occurred between the years 1332-1551.
 229 Introduction to countries west of Nigeria appears to be by dispersal initially to Liberia, followed
 230 by Guinea in the 1700s, followed by Sierra Leone and Mali approximately 100 years later. A
 231 lack of full segment sequences from lineage V limits calculation of divergence from the most
 232 recent common ancestor from lineage IV (approximately 200 years).

235 4 Discussion

236
 237 There are several important aspects of our study and findings. First, we studied a
 238 comprehensive dataset of publicly available full-segment LASV sequences, spanning West
 239 Africa and host species, to inform our understanding of the phylogeny of LASV dispersal.
 240 Second, we identified substantial variability in the origin of available sequences and
 241 completeness of records. Third, we showed strong geographic clustering among lineages
 242 supporting prior hypotheses of radiation from both Nigeria and a subsequent introduction into
 243 Liberia (19). Fourth, the synthesis of available metadata highlights important gaps in currently

244 available data, including spatial bias in the sequencing of samples and suggests this should
245 be used to inform the design of epidemiological programmes going forward.

246

247 Our analyses of 2,298 LASV sequences obtained from GenBank highlights the spatial biases
248 in the availability of sequence data that may limit our understanding of the current and historic
249 dispersal of LASV lineages in West Africa. First, sequence data was typically obtained from
250 three of the eight endemic countries: Nigeria, Guinea and Sierra Leone. We found a strong
251 association between the number of reported human cases and number of available sequences.
252 When stratifying by host species this trend did not remain with rodent derived samples
253 showing no association with the number of human cases indicating important under-sampling
254 in high human cases regions and relatively high sampling in locations with low numbers of
255 human cases. This is potentially an important source of bias when attempting to infer
256 phylogeography within the reservoir host of this zoonotic pathogen. Sequence data from other
257 countries, and more regions within them, across West Africa are required to increase
258 confidence in the timelines of the currently inferred westward expansion. Greater focus needs
259 to be placed on acquiring sequences from the rodent host to understand viral genetic diversity
260 within the primary reservoir species. Comparing rodent derived sequences with those
261 acquired from spillover into human populations may also allow identification of genetic drivers
262 of transmission (36).

263

264 The overrepresentation of data from these three countries has been mapped as relative
265 sequencing effort to identify regions where increased LASV sequencing are required to
266 counteract current sequencing biases. Second, geographic clustering of LASV lineages,
267 suggest isolated events of human-to-rodent transmission and the emergence of LASV dating
268 from 1498 in Nigeria. Similarly, Olayemi *et al.* report evidence of earlier emergence of the virus
269 in humans than in rodents in Nigeria (16). Comparatively limited data from non-human hosts
270 with limited genome coverage, (69/703 sequences encompassed complete genes) produce
271 important uncertainty around the observation of human-to-rodent transmission. Taken
272 together, this data highlight limited surveillance among animal species, necessitating further
273 investments in data acquisition and sharing to accurately define the spatiotemporal expansion
274 of LASV in West Africa.

275

276 The phylogenetic analysis of LASV stratified by host species supports spatial evolution, in
277 addition to intra-host viral evolution (S2 Fig). For instance, LASV sequences from *M.*
278 *erythroleucus* sampled in Nigeria and Guinea clustered within lineages III and IV, respectively.
279 Interestingly, these isolates appear to occur after the emergence of the most recent common
280 ancestor virus circulating among humans and *M. natalensis* in these countries (Table 1),

281 suggesting introduction of LASV into *M. erythroleucus* populations was a consequence of
282 pathogen circulation in human and *M. natalensis* populations. Sequences from *M. natalensis*
283 in Sierra Leone exhibit minimal clustering, and were interspersed with sequences from
284 humans, potentially representing isolated events of pathogen introduction into human
285 populations with spillback into commensal rodent populations (i.e., reverse zoonosis). The
286 most recent common ancestor of LASV sequences from *M. natalensis* in Sierra Leone suggest
287 a later emergence of the virus in this country. Our findings corroborate those of Olayemi et al.,
288 that within Sierra Leone LASV appears to have emerged in human hosts before rodents (16).
289 However, this data must be caveated by the limited information from rodent species in these
290 locations.

291

292 There is a lower coverage of rodent-derived LASV sequences, with those from the primary
293 reservoir *M. natalensis* forming fewer than one-third of all sequences ($n = 642$, 28%), with
294 substantially lower sampling of other possible rodent hosts, including other *Mastomys* species.
295 Rodent sampling has not increased at the same rate as human samples despite increased
296 sampling effort since 2008 (15,22,37). There is substantial heterogeneity in the locations in
297 which rodent and human samples are available. For example, a relatively high number of
298 rodent samples ($n = 429$) have been obtained from Guinea while few human sequences ($n =$
299 20) are available from these locations. The inverse is true of Nigeria where most human
300 derived sequences are obtained ($n = 1,147$) but only 85 rodent sequences are available, and
301 all of these from a single state (Edo, Nigeria). The number of suspected and reported cases
302 was found to be positively but non-linearly associated with the number of available sequences.
303 This is suggestive of a consolidation of research and focus of sampling in areas historically
304 with high numbers of human cases but has led to a paucity of sequences from elsewhere in
305 the endemic region. The limited number of full segment sequences from rodents, from few
306 geographic locations, limits our understanding of viral radiation in rodent hosts, particularly
307 from species which are not considered the primary reservoir, e.g., *H. pamfi*. The most recent
308 common ancestor for the viral sequence obtained from *H. pamfi* is estimated to be in the late
309 1600s, it is therefore possible lineage VI and/or *H. pamfi* as a reservoir of LASV has gone
310 undetected due to lack of sufficient sampling (15).

311

312 Interpreting available LASV sequences is challenging for several reasons. A large proportion
313 of available sequences (70%) have been obtained within Lassa fever research programs,
314 representing spatial ascertainment bias (38–40). In addition to these spatial biases' temporal
315 biases are apparent. Since 2016 there has been a substantial increase in the number of LASV
316 sequences available in NCBI GenBank, reflecting increasing research effort, availability of
317 sequencing platforms and increased data collection during Lassa fever epidemics, such as in

318 the 2018 Nigeria Lassa fever outbreak (41–43). There are notably fewer recorded sequences
319 of LASV from Benin, Togo, and Ghana, suggesting a potential a gap in surveillance and
320 research capacity in these locations or a lack of circulating LASV, despite several reported
321 outbreaks (44–46). Phylogenetic analysis on 60% of our initial dataset, following removal of
322 sequences due to incompleteness or missing geographic and year of collection information (n
323 = 1,045) demonstrated geographic clustering of LASV lineages, supporting prior analyses
324 (14–16,33,44,47–49). Increased data availability from Nigeria following increased LASV
325 surveillance allowed regional analysis of phylogeny for lineages II and III supporting previous
326 findings of expansion of these lineages from North-East Nigeria to the South-West of the
327 country (13,50,51).

328

329 A substantial number ($n = 869$) of the sequences retrieved corresponded to short fragments
330 (< 1 Kb) probably derived from PCR products used for diagnostic purposes rather than for
331 viral genomic surveillance. LASV is a segmented virus, and it was not possible to identify
332 complete genome sequences since both S and L segments are reported separately on the
333 sequence's repository. The molecular clock analyses from L protein indicated an earlier
334 emergence of LASV when compared to S segment analysis (828 and 1350 respectively),
335 potentially because the viral RNA polymerase (L protein) is less affected by selective
336 pressure than the S segment (12,47,52).

337

338 Despite these challenges, this study has synthesised currently available data on LASV
339 sequences to investigate the location and period of sampling to reconstruct the dispersal of
340 viral lineages across the endemic region. Despite the regionalisation of LF being driven by
341 rodent-to-human transmission, there remains scarce LASV genomic data from non-human
342 hosts. We have mapped the locations of relative under sampling to guide targeted efforts to
343 counteract biases in currently available data for both rodent and human derived sequences.
344 Expanded sampling of LASV from animal species within the endemic region will improve our
345 current understanding of LASV evolution and ecology and improve confidence in current
346 estimates of westward expansion of Lassa fever in humans. Further understanding of the
347 viral evolution dynamics of LASV and spatial expansion of current lineages will be vital to
348 ensure adequate diagnostic tools are available to respond to the expected sporadic
349 outbreaks of Lassa Fever across the region.

350

351 **Supplementary material**

352

353 **S1 Data. GenBank accession number of analysed sequences.** This dataset includes
354 available data about host, country, region, year, sequence length, genome segment (L or S)
355 and predicted MRCA.

356

357 **S2 Data. Dataset on confirmed Lassa fever cases.** This presents the number of confirmed
358 cases of Lassa fever reported from countries between 2008 and 2023 at a subnational level
359 that were used to calculate the number of cases per 100,000 people. References for the
360 reports used to produce this dataset are included.

361

362 **S1 Figure. Map of West Africa.** displays a map of West Africa with country names for
363 reference with Fig 1 and Fig 2. Shapefiles for mapping obtained from GADM 4.0.2 (35)

364

365 **S2 Figure. Time-calibrated phylogeny for both the small segment (S) and large segment
366 (L) from included LASV sequences.**

367

368 **Author contributions**

369 Conceptualisation: DS and LBA; Methodology: HF, DS, DA and LBA; Formal Analyses: HF,
370 DS and LBA; Investigation: HF, DS and LBA; Supervision: LBA; Data Curation: HF and DS;
371 Writing – original draft preparation: HF, DS, LBA; Writing – Review and Editing: IH, LE, NH,
372 RA, RK, FN, DA, AZ and TMcH; Funding acquisition: AZ and FN.

373

374 **Data availability and reproducibility**

375 All data used in these analyses are publicly available from GenBank. The accession numbers
376 of records used are available as supplementary material. Code to reproduce the metadata
377 analyses are available as an archived Git release on Zenodo
378 (<https://doi.org/10.5281/zenodo.6340162>)

379

380 **Conflict of interests**

381 The authors declare no conflict of interests

382

383 **Acknowledgements:**

384 Linzy Elton, Timothy D McHugh, Francine Ntoumi, and Alimuddin Zumla acknowledge support
385 from EDCTP-Central Africa and East African Clinical Research Networks (CANTAM-3,
386 EACCR-3). Sir Zumla is an NIHR Senior Investigator, a Mahathir Science Award, Sir Patrick
387 Manson Medal and EU-EDCTP Pascoal Mocumbi Prize laureate. Liã Bárbara Arruda, David
388 Simons, Rashid Ansumana, Linzy Elton, Najmul Haider, Isobella Honeyborne, Danny Asogun,
389 Timothy D McHugh, Francine Ntoumi, Alimuddin Zumla and Richard Kock acknowledge

390 support from the Pan-African Network for Rapid Research, Response and Preparedness for
391 Infectious Diseases Epidemics – PANDORA-ID-NET, funded through the European and
392 Developing Countries Clinical Trials Partnership (EDCTP) (grant number RIA2016E-1609).
393 David Simons is supported by a PhD studentship from the UK Biotechnology and Biological
394 Sciences Research Council (BB/M009513/1).

395

396 **References**

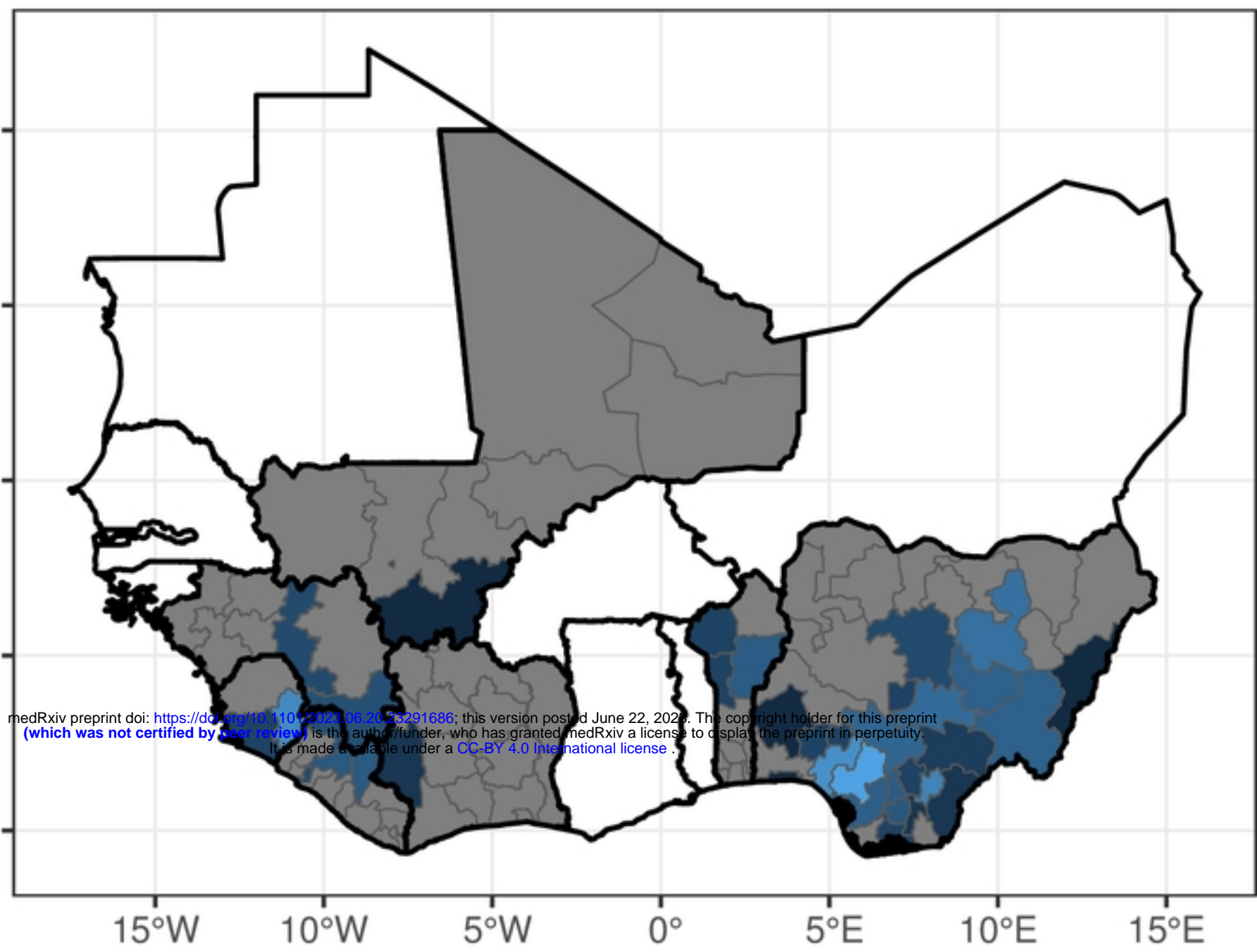
- 397 1. Asogun DA, Gunther S, Akpede GO, Ihekweazu C, Zumla A. Lassa Fever: Epidemiology, Clinical
398 Features, Diagnosis, Management and Prevention. [Review]. *Infectious Disease Clinics of North*
399 *America*. 2019;33(4):933–51.
- 400 2. Basinski AJ, Fichet-Calvet E, Sjodin AR, Varrelman TJ, Remien CH, Layman NC, et al. Bridging the
401 gap: Using reservoir ecology and human serosurveys to estimate Lassa virus spillover in West
402 Africa. *Wesolowski A, editor. PLoS Comput Biol*. 2021 Mar 3;17(3):e1008811.
- 403 3. World Health Organisation. Lassa fever [Internet]. 2022 [cited 2022 Feb 22]. Available from:
404 https://www.who.int/health-topics/lassa-fever#tab=tab_1
- 405 4. Gunther S, Lenz O. Lassa virus. [Review] [323 refs]. *Critical Reviews in Clinical Laboratory*
406 *Sciences*. 2004;41(4):339–90.
- 407 5. Hallam SJ, Koma T, Maruyama J, Paessler S. Review of Mammarenavirus Biology and Replication.
408 *Front Microbiol* [Internet]. 2018 [cited 2020 Oct 21];9. Available from:
409 <https://www.frontiersin.org/articles/10.3389/fmicb.2018.01751/full>
- 410 6. Welch SR, Scholte FEM, Albariño CG, Kainulainen MH, Coleman-McCray JD, Guerrero LW, et al.
411 The S Genome Segment Is Sufficient to Maintain Pathogenicity in Intra-Clade Lassa Virus
412 Reassortants in a Guinea Pig Model. *Frontiers in Cellular and Infection Microbiology* [Internet].
413 2018 [cited 2022 Feb 3];8. Available from:
414 <https://www.frontiersin.org/article/10.3389/fcimb.2018.00240>
- 415 7. Simons D. Lassa fever cases suffer from severe underreporting based on reported fatalities.
416 *International Health*. 2022;
- 417 8. McCormick JB, Webb PA, Krebs JW, Johnson KM, Smith ES. A prospective study of the
418 epidemiology and ecology of Lassa fever. *J Infect Dis*. 1987;155(3):437–44.
- 419 9. Takah NF, Brangel P, Shrestha P, Peeling R. Sensitivity and specificity of diagnostic tests for Lassa
420 fever: a systematic review. *BMC Infectious Diseases*. 2019 Jul 19;19(1):647.
- 421 10. Nnaji ND, Onyeaka H, Reuben RC, Uwishema O, Olovo CV, Anyogu A. The deuce-ace of Lassa
422 Fever, Ebola virus disease and COVID-19 simultaneous infections and epidemics in West Africa:
423 clinical and public health implications. *Tropical Medicine and Health*. 2021 Dec 30;49(1):102.
- 424 11. Ashcroft JW, Olayinka A, Ndodo N, Lewandowski K, Curran MD, Nwafor CD, et al. Pathogens that
425 Cause Illness Clinically Indistinguishable from Lassa Fever, Nigeria, 2018. *Emerging Infectious*
426 *Diseases*. 2022;28(5):994–7.
- 427 12. Andersen KG, Shapiro BJ, Matranga CB, Sealfon R, Lin AE, Moses LM, et al. Clinical Sequencing
428 Uncovers Origins and Evolution of Lassa Virus. *Cell*. 2015;

- 429 13. Bowen MD, Rollin PE, Ksiazek TG, Hustad HL, Bausch DG, Demby AH, et al. Genetic Diversity
430 among Lassa Virus Strains. *Journal of Virology*. 2000;74(15):6992–7004.
- 431 14. Manning JT, Forrester N, Paessler S. Lassa virus isolates from Mali and the Ivory Coast represent
432 an emerging fifth lineage. *Frontiers in Microbiology*. 2015;
- 433 15. Olayemi A, Cadar D, Magassouba N, Obadare A, Kourouma F, Oyeyiola A, et al. New Hosts of The
434 Lassa Virus. *Scientific Reports*. 2016;
- 435 16. Olayemi A, Adesina AS, Strecker T, Magassouba N, Fichet-Calvet E. Determining ancestry
436 between rodent-and human-derived virus sequences in endemic foci: Towards a more integral
437 molecular epidemiology of lassa fever within West Africa. *Biology*. 2020;
- 438 17. Whitmer SLM, Strecker T, Cadar D, Dienes HP, Faber K, Patel K, et al. New lineage of lassa virus,
439 Togo, 2016. *Emerging Infectious Diseases*. 2018;24(3):599–602.
- 440 18. Okoro OA, Bamgboye E, Dan-Nwafor C, Umeokonkwo C, Ilori E, Yashe R, et al. Descriptive
441 epidemiology of Lassa fever in Nigeria, 2012-2017. *Pan Afr Med J*. 2020 Sep 3;37:15.
- 442 19. Klitting R, Kafetzopoulou LE, Thiery W, Dudas G, Gryseels S, Kotamarthi A, et al. Predicting the
443 evolution of Lassa Virus endemic area and population at risk over the next decades [Internet].
444 *Microbiology*; 2021 Sep [cited 2022 Feb 3]. Available from:
445 <http://biorxiv.org/lookup/doi/10.1101/2021.09.22.461380>
- 446 20. Bangura U, Buanie J, Lamin J, Davis C, Bongo GN, Dawson M, et al. Lassa Virus Circulation in
447 Small Mammal Populations in Bo District, Sierra Leone [Internet]. Vol. 10, *BIOLOGY-BASEL. ST*
448 *ALBAN-ANLAGE 66, CH-4052 BASEL, SWITZERLAND: MDPI*; 2021. Available from:
449 <https://doi.org/10.3390/biology10010028>
- 450 21. Forni D, Sironi M. Population Structure of Lassa Mammarenavirus in West Africa. *Viruses*.
451 2020;12(4):437.
- 452 22. Lecompte E, Fichet-Calvet E, Daffis S, Koulémou K, Sylla O, Kourouma F, et al. *Mastomys*
453 *natalensis* and Lassa fever, West Africa. *Emerging Infectious Diseases*. 2006;
- 454 23. Wulff H, Fabiyi A, Monath TP. Recent isolations of Lassa virus from Nigerian rodents. *Bull World*
455 *Health Organ*. 1975;52(4–6):609–13.
- 456 24. Yadouleton A, Agolinou A, Kourouma F, Saizonou R, Pahlmann M, Bedié SK, et al. Lassa virus in
457 pygmy mice, Benin, 2016-2017. *Emerging Infectious Diseases*. 2019;
- 458 25. Oti VB. A Reemerging Lassa Virus: Aspects of Its Structure, Replication, Pathogenicity and
459 Diagnosis. In: Alfonso J. Rodriguez-Morales, editor. *Current Topics in Tropical Emerging Diseases*
460 *and Travel Medicine*. BoD – Books on Demand; 2018.
- 461 26. Lo Iacono G, Cunningham AA, Fichet-Calvet E, Garry RF, Grant DS, Khan SH, et al. Using
462 Modelling to Disentangle the Relative Contributions of Zoonotic and Anthroponotic
463 Transmission: The Case of Lassa Fever. *PLoS Neglected Tropical Diseases*. 2015;
- 464 27. National Center for Biotechnology Information. National Center for Biotechnology Information
465 [Internet]. 2022 [cited 2022 Feb 3]. Available from: <https://www.ncbi.nlm.nih.gov/>
- 466 28. Becker G, Lawrence M. *genbankr: Parsing GenBank files into semantically useful objects*. 2021.

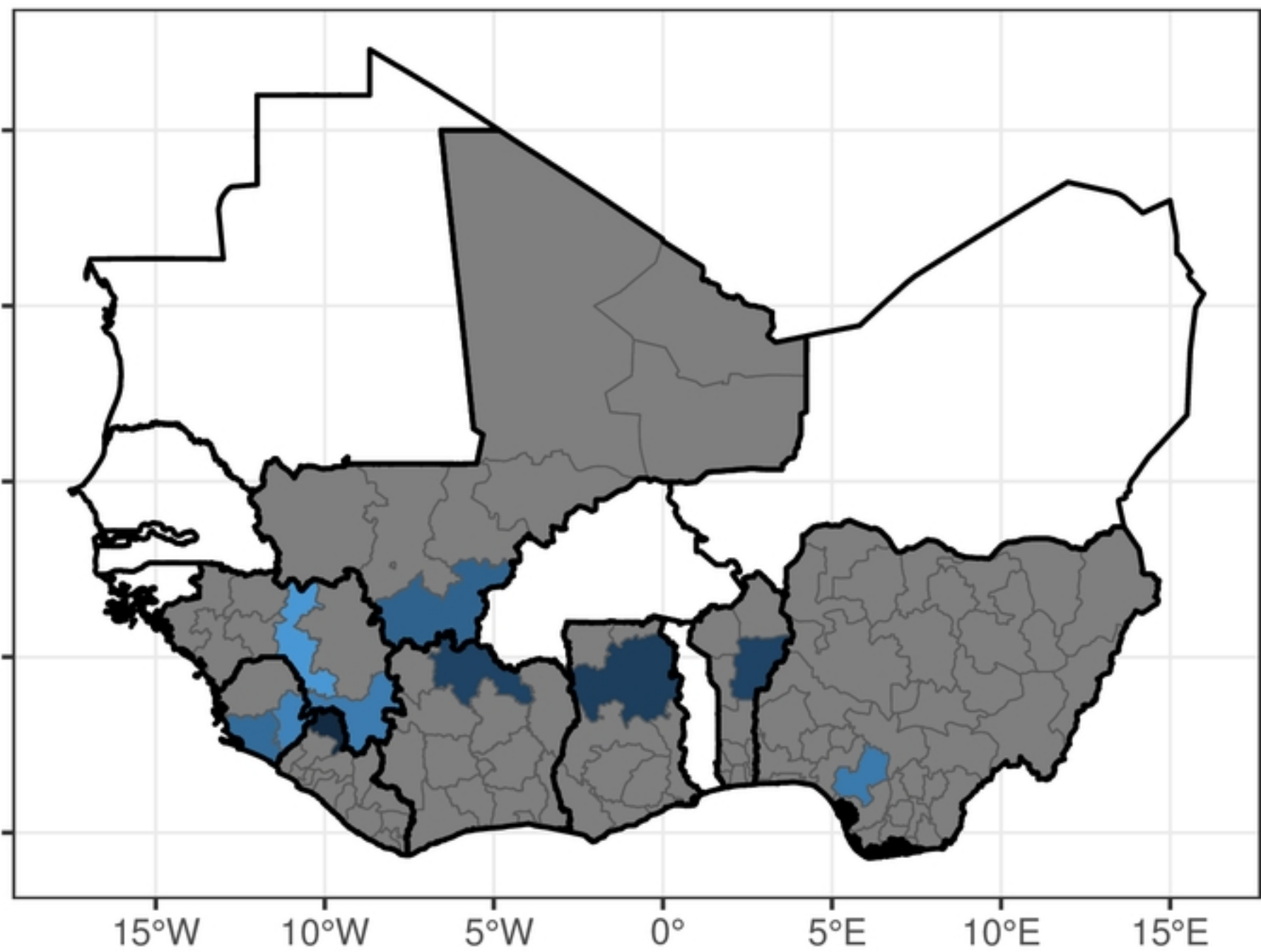
- 467 29. R Core Team. R: A Language and Environment for Statistical Computing [Internet]. Vienna,
468 Austria: R Foundation for Statistical Computing; 2021. Available from: [https://www.R-](https://www.R-project.org/)
469 [project.org/](https://www.R-project.org/)
- 470 30. Kahle D, Wickham H. ggmap: Spatial Visualization with ggplot2. *The R Journal*. 2013;5(1):144–61.
- 471 31. Wood SN. *Generalized Additive Models: An Introduction with R*. 2nd ed. Chapman and Hall/CRC;
472 2017.
- 473 32. Suchard MA, Lemey P, Baele G, Ayres DL, Drummond AJ, Rambaut A. Bayesian phylogenetic and
474 phylodynamic data integration using BEAST 1.10. *Virus Evolution*. 2018;
- 475 33. Olayemi A, Fichet-Calvet E. Systematics, ecology, and host switching: Attributes affecting
476 emergence of the Lassa virus in rodents across western Africa. *Viruses*. 2020.
- 477 34. Rambaut A, Drummond AJ, Xie D, Baele G, Suchard MA. Posterior summarization in Bayesian
478 phylogenetics using Tracer 1.7. *Systematic Biology*. 2018;
- 479 35. Database of Global Administrative Areas. GADM [Internet]. 2022 [cited 2021 Apr 25]. Available
480 from: <https://gadm.org/index.html>
- 481 36. Whitlock AOB, Bird BH, Ghersi B, Davison AJ, Hughes J, Nichols J, et al. Identifying the genetic
482 basis of viral spillover using Lassa virus as a test case. *R Soc Open Sci*. 2023 Mar 22;10(3):221503.
- 483 37. Lecompte E, Brouat C, Duplantier JM, Galan M, Granjon L, Loiseau A, et al. Molecular
484 identification of four cryptic species of *Mastomys* (Rodentia, Murinae). *Biochemical Systematics*
485 *and Ecology*. 2005;
- 486 38. Townsend Peterson A, Moses LM, Bausch DG. Mapping transmission risk of lassa fever in West
487 Africa: The importance of quality control, sampling bias, and error weighting. *PLoS ONE*. 2014;
- 488 39. Ehichioya DU, Hass M, Ölschläger S, Becker-Ziaja B, Onyebuchi Chukwu CO, Coker J, et al. Lassa
489 fever, Nigeria, 2005-2008. *Emerging Infectious Diseases*. 2010.
- 490 40. Khan SH, Goba A, Chu M, Roth C, Healing T, Marx A, et al. New opportunities for field research
491 on the pathogenesis and treatment of Lassa fever. *Antiviral Research*. 2008;
- 492 41. Maxmen A. Deadly Lassa-fever outbreak tests Nigeria's revamped health agency. *Nature*.
493 2018;555(7697):421–2.
- 494 42. Siddle KJ, Eromon P, Barnes KG, Mehta S, Oguzie JU, Odia I, et al. Genomic Analysis of Lassa
495 Virus during an Increase in Cases in Nigeria in 2018. *New England Journal of Medicine*. 2018 Nov
496 1;379(18):1745–53.
- 497 43. Ilori EA, Frank C, Dan-Nwafor CC, Ipadeola O, Krings A, Ukponu W, et al. Increase in Lassa Fever
498 Cases in Nigeria, January–March 2018. *Emerging Infectious Diseases* [Internet]. 2019 May [cited
499 2020 Oct 21];25(5). Available from: 10.3201/eid2505.181247
- 500 44. Yadouleton A, Picard C, Rieger T, Loko F, Cadar D, Kouthon EC, et al. Lassa fever in Benin:
501 description of the 2014 and 2016 epidemics and genetic characterization of a new Lassa virus.
502 *Emerging Microbes & Infections*. 2020;1–23.

- 503 45. World Health Organisation. Lassa Fever – Togo [Internet]. 2022 [cited 2022 Nov 24]. Available
504 from: <https://www.who.int/emergencies/disease-outbreak-news/item/2022-DON362>
- 505 46. Ghana Health Services. Lassa Fever Press Release Ghana 2023. 2023 Apr 19 [cited 2023 Apr 19];
506 Available from: <https://osf.io/ft2gy/>
- 507 47. Ibukun FI. Inter-lineage variation of lassa virus glycoprotein epitopes: A challenge to lassa virus
508 vaccine development. *Viruses*. 2020.
- 509 48. Lalis A, Leblois R, Lecompte E, Denys C, ter Meulen J, Wirth T. The Impact of Human Conflict on
510 the Genetics of *Mastomys natalensis* and Lassa Virus in West Africa. *PLoS ONE*. 2013;7(5).
- 511 49. Wiley MR, Fakoli L, Letizia AG, Welch SR, Ladner JT, Prieto K, et al. Lassa virus circulating in
512 Liberia: a retrospective genomic characterisation. *The Lancet Infectious Diseases*. 2019;
- 513 50. Ehichioya DU, Hass M, Becker-Ziaja B, Ehimuan J, Asogun DA, Fichet-Calvet E, et al. Current
514 molecular epidemiology of Lassa virus in Nigeria. *Journal of Clinical Microbiology*. 2011;
- 515 51. Naidoo D, Ihekweazu C. Nigeria's efforts to strengthen laboratory diagnostics – Why access to
516 reliable and affordable diagnostics is key to building resilient laboratory systems. *African Journal*
517 *of Laboratory Medicine* [Internet]. 2020 Aug 26 [cited 2020 Oct 21];9(2). Available from:
518 [10.4102/ajlm.v9i2.1019](https://doi.org/10.4102/ajlm.v9i2.1019)
- 519 52. Hastie KM, Saphire EO. Lassa virus glycoprotein: stopping a moving target. *Current Opinion in*
520 *Virology*. 2018.
- 521

Human



Rodent



Number of sequences (log₁₀)

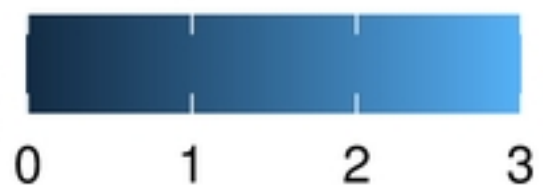


Figure 1

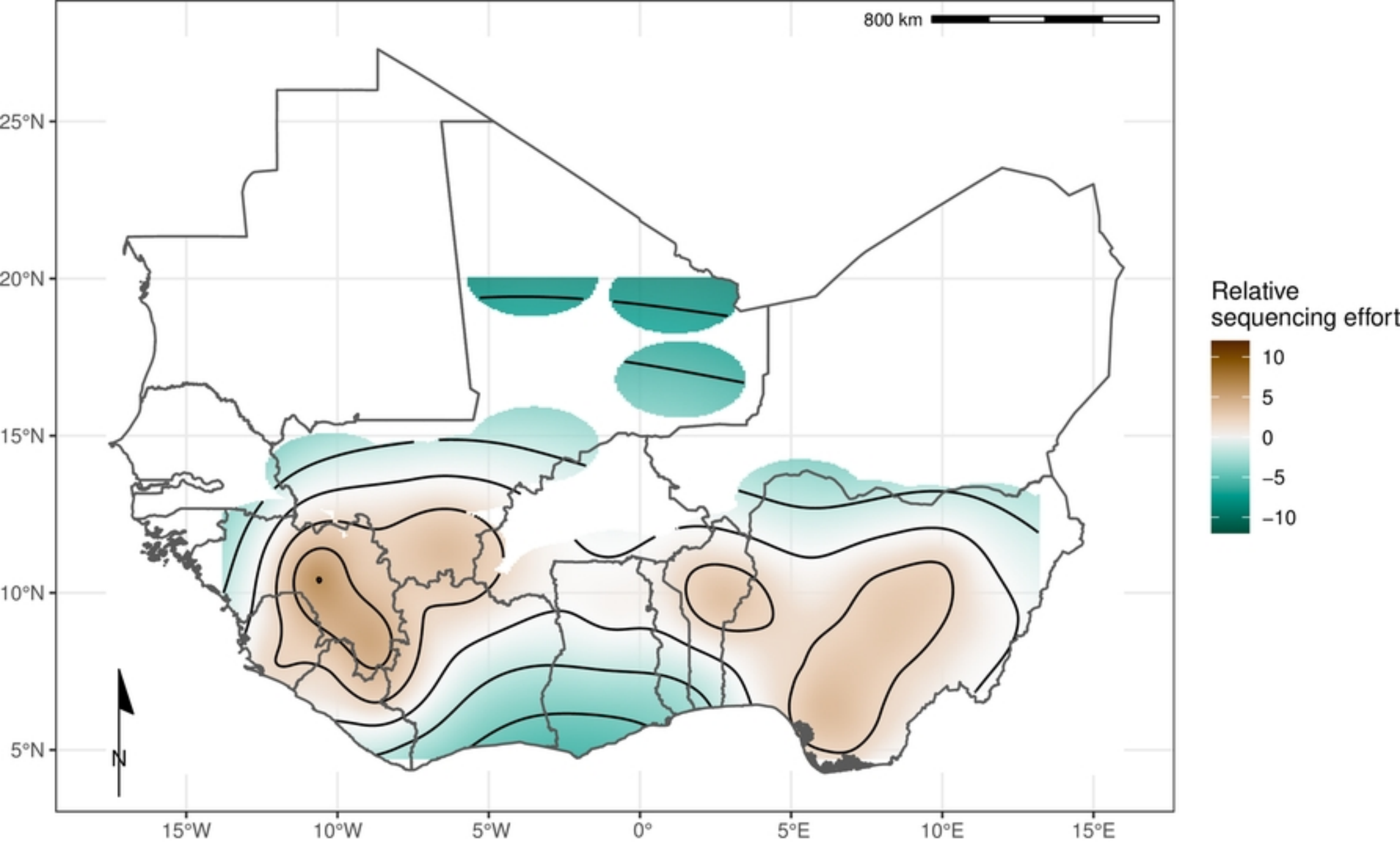


Figure 2