

1 **Genomic epidemiology and evolutionary dynamics of respiratory syncytial virus group B in**

2 **Kilifi, Kenya, 2015-17**

3

4 Everlyn Kamau ^{1,*,#}, James R. Otieno ¹, Nickson Murunga ¹, John W. Oketch ¹, Joyce M. Ngoi

5 ¹, Zaydah R. de Laurent ¹, Anthony Mwema ¹, Joyce U. Nyiro ¹, Charles N. Agoti ^{1,2}, D. James

6 Nokes ^{1,3}

7

8 ¹ Epidemiology and Demography Department, KEMRI–Wellcome Trust Research Programme, Kilifi, Kenya

9 ² Pwani University, School of Health and Human Sciences, Kilifi, Kenya.

10 ³ School of Life Sciences and Zeeman Institute (SBIDER), University of Warwick, Coventry, UK

11 * Correspondence: everlyn.kamau@ndm.ox.ac.uk

12 # Current address: Nuffield Department of Medicine, University of Oxford, Oxford, UK

13

14

15 **Abstract**

16 Respiratory syncytial virus (RSV) circulates worldwide and is a leading cause of acute
17 respiratory illness in young children. There is paucity of genomic data from purposively
18 sampled populations by which to investigate evolutionary dynamics and transmission
19 patterns of RSV. Here we present an analysis of 295 RSV group B genomes from Kilifi,
20 coastal Kenya, sampled from individuals seeking outpatient care in 9 health facilities across
21 a defined geographical area (890 km²), over 2 RSV epidemics between 2015 and 2017. RSVB
22 diversity was characterized by multiple viral introductions into the area and co-circulation of
23 distinct genetic groups or clusters, which transmitted and diversified locally but with varying
24 frequency. Bayesian analyses indicated a strong spatially and temporally structured viral
25 population suggesting extensive within-epidemic virus transmission. Phylogeographic
26 analysis provided a strong support for epidemiological linkage from one central health
27 facility to other facilities. Increase in relative diversity paralleled increase in seasonal viral
28 incidence. Importantly, we identified a cluster of viruses (n=91) that emerged in the
29 2016/17 epidemic, carrying distinct amino-acid signatures including a novel non-
30 synonymous change (K68Q) in antigenic site Ø in the Fusion gene. A different non-
31 synonymous change K68N was recently associated with escape from a potent neutralizing
32 monoclonal antibody (MEDI8897). RSVB diversity was additionally marked by signature non-
33 synonymous substitutions that were unique to particular genomic clusters, some of which
34 were under diversifying selection. Our findings provide insights into recent evolutionary and
35 epidemiologic behaviors of RSV group B, and highlight possible emergence of a novel
36 antigenic variant, which has implications on current prophylactic development strategies.

37 Introduction

38 Respiratory syncytial virus (RSV) is the most common cause of acute lower
39 respiratory tract infection in children <5 years of age worldwide, with an estimated
40 associated mortality of up to 199,000 deaths per year mostly in developing countries (1,2).
41 RSV is also an important cause of community-acquired pneumonia among hospitalized
42 adults of all ages (3). RSV is a member of the family *Pneumoviridae*, subfamily
43 *Orthopneumovirus*, has an enveloped, non-segmented, single-stranded, negative sense RNA
44 genome of approximately 15,000 nucleotides encoding 11 proteins: NS1, NS2, N, P, M, SH,
45 G, F, M2-1, M2-2, and L (4). RSV clinical isolates are classified into two groups (RSVA and
46 RSVB) based on antigenic and genetic variability (5). Distinct genotypes of RSV have been
47 shown to circulate locally and globally (6-9) and strains may vary from location to location in
48 any given season with viruses identified in one location being similar to those from vastly
49 different geographic locations identified in different years (10), suggestive of rapid global
50 transmission. The available therapeutic modalities for RSV are chiefly supportive, and
51 prophylactic treatment with RSV-specific neutralizing antibody is effective in reducing RSV
52 morbidity in infants, but its use is currently limited to high-risk populations in high-resource
53 settings (11). There is no licensed RSV vaccine for routine use in immunization, however,
54 prophylactic vaccine candidates and monoclonal antibodies (mAbs) are now in advanced
55 clinical trials (12).

56 We have previously characterised the community dynamics of RSV in the coastal
57 region of Kenya, using sequences of the highly variable G glycoprotein gene (encoding the
58 attachment protein) for RSV groups A and B and using whole genomes of RSVA genotype
59 ON1, almost exclusively from samples from pneumonia patients admitted to the Kilifi
60 County Hospital (13-15). From these studies, RSV displays high genetic diversity of locally

61 circulating strains, within and between consecutive epidemics. Furthermore, recurrent RSV
62 epidemics in Kilifi are depicted by sequential replacement of genotypes, over the long term,
63 and high turnover of variants within genotypes in the short term (13,14). In the study
64 reported here, samples arise from a design aimed to limit temporal, age-related, illness
65 severity, geographical, and health care access bias. Recruitment was carried throughout a
66 study location, from representative health facilities, simultaneously, and of patients of any
67 age with mild acute respiratory symptoms (16).

68 Phylogenetic and phylogeographic methods have been used increasingly to study
69 molecular epidemiology and evolutionary dynamics of virus populations e.g. Ebola, Zika, and
70 Influenza (17-21). However, despite the importance of RSV to pneumonia hospitalisation
71 and mortality among children (1), equivalent genome-scale studies to examine RSV
72 transmission patterns and evolutionary dynamics within a community setting, are still few
73 and remain a major gap in our understanding of what viral factors shape epidemiologic and
74 evolutionary processes RSV at the local level. While most molecular studies on RSV have
75 focused on the G glycoprotein gene because of its high genetic diversity and utility as a
76 phylogenetic marker, genome-wide genetic signatures in RSV genome additionally inform
77 on diversity and the adaptive mechanisms of RSV viruses following introduction into
78 communities (15).

79 Here, we applied a combination of molecular clock, coalescent and discrete diffusion
80 phylogeographic models to whole genome sequences to measure genomic diversity and
81 deduce spatial and temporal circulation of RSVB in rural Kilifi, coastal Kenya. We sought to
82 characterize introductions, transmission and spread of RSVB from samples collected
83 through outpatient surveillance of respiratory viruses, analogous to studying community
84 epidemics. We present estimates of evolutionary parameters such as the genomic rate of

85 evolution and relative genetic diversity and describe spatial and temporal clustering
86 patterns of viruses within each epidemic as demonstrated by phylogenetically distinct
87 clades. In particular, we determine possible selection-driven emergence of a novel RSVB
88 variant carrying distinct amino acid signatures. Our study significantly increases the number
89 of publicly available complete RSVB genomes, which will enable further studies of RSV
90 evolution.

91

92 **Materials and methods**

93

94 **Study site and patient recruitment.** The annual RSV epidemics in Kilifi, coastal
95 Kenya, are seasonal starting from November through May, with a peak around January and
96 an average duration of 18 weeks (22). The study from which the samples arise and used in
97 the present report was done from December 2015 to March 2018, a period covering three
98 RSV seasonal epidemics (2015/16, 2016/17 and 2017/18), and was carried out within the
99 Kilifi Health and Demographic Surveillance System (KHDSS) (23). The study was conducted
100 to document the community-wide burden of respiratory virus infections. Nine out of the 24
101 public outpatient health facilities in KHDSS were purposively selected (**Figure 1(A)**) –
102 Matsangoni (MAT), Ngerenya (NGE), Mtondia (MTO), Sokoke (SOK), Mavueni (MAV),
103 Jaribuni (JAR), Chasimba (CHA), Pingilikani (PIN) and Junju (JUN) - to provide a broad
104 representation across the geographical region, covering major road networks into the
105 location and variation in population density (16). Participant recruitment and specimen
106 collection was integrated within the routine patient care and led by a resident clinician or
107 nurse as detailed in (16). Patients of any age presenting with one or more ARI symptoms of
108 cough, sneezing, nasal congestion, difficulty breathing, or increased respiratory rate for age

109 were eligible. Written individual informed consent was sought from adult patients and
110 parents/guardians of patients below 18 years.

111 The study was approved by the KEMRI-Scientific and Ethical Review Unit (SERU#
112 3103) and the University of Warwick Biomedical and Scientific Research Ethics Committee
113 (BSREC# REGO-2015-6102).

114

115 **Sample collection and laboratory testing.** A nasopharyngeal swab (NPS) was
116 collected from each participant and stored in universal virus transport media (Copan
117 Diagnostics, USA). RNA was extracted from samples by Qiacube HT using an RNeasy
118 extraction kit (Qiagen, Germany) and screened for a range of respiratory viruses including
119 RSV A and RSV B using a multiplex real-time PCR assay system (24,25). Samples with real time
120 RT-PCR cycle threshold (Ct) of <35.0 were defined as positive for the target virus.

121

122 **RSV B whole genome amplification and sequencing.** Whole genome amplification
123 and sequencing was attempted for all RSV B positive samples with PCR cycle threshold (Ct)
124 value < 35.0 collected between mid-December 2015 to end of May 2017. No RSV B was
125 detected from June 2017 to the end of the study in March 2018 (i.e. 2017/18 epidemic was
126 RSV A only). Reverse transcription of RNA molecules and PCR amplification were performed
127 with a six-amplicon, six-reaction strategy presented in detail in (26). Briefly, extracted RNA
128 was converted to cDNA using Superscript III First Strand Synthesis kit (Invitrogen) and
129 forward primers, followed by PCR using Phusion High-Fidelity DNA polymerase (NEB).
130 Amplification success was confirmed by observing the expected PCR product size (2300–
131 2500 bp) on 0.6% agarose gels. Amplicons were fragmented, tagged to adapters, and
132 indexed using the Nextera XT (Illumina, San Diego, CA, USA) library prep kit as per

133 manufacturer's instructions. Size distribution of the barcoded libraries was assessed by
134 Agilent's 2100 Bioanalyzer. Pooled and normalized libraries in batches of seventy or eighty
135 were sequenced on Illumina MiSeq system using 2 x 250 bp or 2 x 300 bp paired-end (PE)
136 sequencing at the KEMRI-Wellcome Trust Research Program laboratories.

137

138 **Short read data assembly.** Methods for quality check of the sequence reads,
139 depletion of human reads, generation of consensus genome assemblies and calculation of
140 coverage depth, were as described in (15). Briefly, quality check of the sequence reads was
141 done using FastQC. Consensus genome assemblies were generated using viral-ngs pipeline
142 v1.19.0 (Broad Institute). The raw reads were depleted of human reads by mapping onto the
143 human reference genome hg19 using bowtie2 (27), and samtools (28) used to filter, sort,
144 and recover the unmapped reads. The raw reads were mapped onto individual consensus
145 assemblies with bowtie2, then samtools used to sort and index the aligned .bam files, and
146 BEDtools (29) used to generate the coverage depth statistics.

147

148 **Sequence data compilation.** To contextualize the diversity of RSVB in Kilifi in 2015 to
149 2017, the new Kilifi genomes were appended to a global dataset of other RSVB genomes
150 available in Genbank (>14000 nucleotides (nt)), retrieved on 29 September 2019
151 (**Supplementary Table 1**). Only global sequences belonging to the RSVB BA genotype were
152 considered, since the Kilifi sequences were of the BA genotype. The global dataset
153 comprised sequences of viruses collected between 2012 and 2016 from UK, US, Russia,
154 Philippines, Nicaragua, Jordan and India. Sequences of viruses sampled before 2012 and
155 published without date or location of sampling were excluded.

156

157 **Phylogenetic and phylogeographic analyses.** Sequence alignment was done using
158 MAFFT v.7.221 (30) using the parameters ‘–localpair –maxiterate 1000’. To determine the
159 degree of temporal signal of divergence in the Kilifi RSVB genomes, a maximum likelihood
160 (ML) tree was estimated using RAxML (31) under general time-reversible (GTR) nucleotide
161 substitution model with gamma-distributed among-site rate heterogeneity, and examined
162 by exploratory linear regression analysis in Tempest v1.5.1 (32). jModelTest (33) was used to
163 determine best-fit substitution model. Root-to-tip divergence was plotted as a function of
164 sampling date (day-month-year). Branch support for ML trees was evaluated using 1000
165 bootstrap replicates.

166 Bayesian molecular clock phylogenies and discrete trait phylogeographic analyses
167 were done using BEAST v1.10 (34) for 300 million MCMC steps, sampling parameters and
168 trees every 5000 steps. HKY nucleotide substitution model with a gamma-distributed rate
169 variation among sites was used. Path-sampling and stepping-stone marginal likelihood
170 estimator (MLE) models (35) were used to estimate the most probable combination of
171 uncorrelated lognormal relaxed molecular clock and coalescent models. The best fitting
172 combination was the uncorrelated log-normal distributed relaxed molecular clock and
173 Skygrid’s Gaussian Markov random field (GMRF) tree prior. A non-informative continuous
174 time Markov chain reference prior was used on the molecular clock. Mixing and
175 convergence of the MCMC sampler in the posterior target distribution was evaluated using
176 Tracer v1.6 (<http://beast.bio.ed.ac.uk/Tracer>), and a maximum clade credibility (MCC) trees
177 summarized with TreeAnnotator after removal of 10% burn-in. MCC trees were visualized
178 with FigTree 1.4.4 (<http://tree.bio.ed.ac.uk/software/figtree/>).

179 To describe and quantify viral movement in Kilifi, we performed a discrete trait
180 diffusion analysis, using sampling location as a discrete trait. We used a Bayesian stochastic

181 search variable selection (BSSVS) procedure (36) and compared reversible and non-
182 reversible (asymmetric) discrete diffusion phylogeographic models to estimate the most
183 relevant transition pathways within the study locations. The standard implementation of
184 the discrete asymmetrical phylogeographic approach assigns the same prior probability
185 distribution to all of the transition rates in the continuous time markov chain (CTMC) (37).
186 On the other hand, the reversible CTMC model allows more balanced transitions in the
187 phylogeny or free lineage exchange between location states without preconditioned
188 directionality but may have a poor fit for locations that could have unidirectional links in
189 reality e.g. spatially expanding epidemics (36). Statistical support for migratory events was
190 measured using Bayes factors (BF) and summarized using Spread3 software (38) after
191 discarding 10% burn-in. We assumed that a $BF > 6$ is strong evidence for a well-supported
192 viral pathway between two locations.

193

194 **Molecular evolution and adaptation.** Gene-specific ratio of nonsynonymous to
195 synonymous substitutions per site (dN/dS) were estimated for codon-based alignments
196 using the single likelihood ancestor counting (SLAC) method available at the Datamonkey
197 webserver (39). We also investigated episodic positive or diversifying selection using mixed-
198 effects model of evolution (MEME) and Fast, Unconstrained Bayesian AppRoximation
199 (FUBAR) approaches also available in Datamonkey. MEME aims to detect sites evolving
200 under positive selection in a proportion of branches, while FUBAR uses a Bayesian approach
201 and assumes that the selection pressure is constant along the entire phylogeny.

202

203 **Phylogenomic clusters.** We defined a phylogenomic cluster as genetically linked viral
204 sequences (≥ 2) that were more closely related than any randomly selected sequences

205 within genetic distance threshold of $<d$. Maximum likelihood evolutionary distances, an
206 estimate of the average number of changes per base pair, were calculated for each pair of
207 consensus genomes using IQ-TREE (40) to measure genetic similarity or dissimilarity
208 between pairs of taxa. ML distances, depicting genetic distance between sequences in a
209 phylogenetic tree measured in nucleotide substitutions per site, were selected over the
210 conventional nucleotide p -distances in order to correct for multiple nucleotide
211 substitutions.

212 The distribution histogram plots of pairwise genetic distances were generated in R,
213 and a distance cutoff value was determined as the least frequent value (d) between the first
214 and second peaks in the histogram. Median values of pairwise ML evolutionary distance
215 distributions of the identified phylogenomic cluster were checked if they were below t -
216 percentile of the overall distance distribution. The t -percentile corresponded to the d
217 distance threshold. In addition, a phylogenomic cluster was only considered eligible if it had
218 reliability of $\geq 90\%$ phylogenetic bootstrap support or Bayesian posterior probability. Each
219 inferred phylogenomic cluster was evaluated for size, temporal and genetic characteristics.

220

221 **Phylogeny-trait association analysis.** We used the Bayesian Tip-association
222 Significance (BaTS) software (41), which uses the posterior sets of trees from Bayesian
223 MCMC analysis, to measure the degree of association between sampling location and the
224 phylogeny and are inversely correlated with the degree of association. Each sequence was
225 assigned a location trait reflecting its origin of sampling and the first 10% of tree states were
226 removed as burn-in. The overall statistical significance was determined by estimating the
227 parsimony score (PS) and association index (AI) metrics, where the null hypothesis was that
228 clustering by location was not more than expected by chance alone. In addition, the

229 maximum clade size (MC) metric was used to compare the strength of clustering at each
230 location by calculating the expected (null) and the observed mean clade size from each
231 study location. A significance level of 0.05 was used in all cases. The PS, AI and MC statistics
232 were computed for a null distribution with 1000 replicates.

233

234 **Sequence data availability.** The final sequencing reads are available in the NCBI
235 BioProject database under the study accession PRJNA562116 and the genomes generated in
236 this study are available in GenBank under accession numbers MN365302 to MN365600.

237

238 **Results**

239 **RSV occurrence in Kilifi, 2015 to 2017.** Between December 2015 and June 2017, a
240 total of 8127 nasopharyngeal swab samples were tested for RSV. RSV was detected in 503
241 (6.19%) samples (Ct <35): 95 (18.9%) were RSV group A (RSVA) and 408 (81.1%) were RSV
242 group B (RSVB). The frequency and pattern of occurrence of RSVB for each health facility are
243 shown in **Figure 1(B)**. The mean (SD) Ct value was 26.9 (4.16) and 25.0 (4.16) for RSVA and
244 RSVB, respectively. The proportion of RSV positive individuals differed by age (p value
245 <0.001), study location (p value = 0.003) but not by gender (p value = 0.078) (**Table 1**). The
246 distribution of RSV viral load (equated to rRT-PCR Ct value) did not differ by outpatient
247 health facility (**Supplementary Figure 1(A)**). The median age of RSV positive individuals was
248 20 months (interquartile range (IQR), 8-43 months), 81.7% (411/503) were aged 5 years or
249 younger, and 272 (54.1%) of the cases were female (**Table 1**). RSV prevalence across
250 different age groups by gender is shown in **Supplementary Figure 2**.

251 RSV co-infection with other respiratory viruses was observed in 41/503 (8.2%)
252 samples: rhinovirus (n =18) and adenovirus (n=11) were co-detected most frequently with

253 RSV. Comparative RSV incidence data was obtained from a contemporaneous respiratory
254 virus surveillance study that recruits patients aged <59 months admitted with acute lower
255 respiratory tract infection at the pediatric ward of the Kilifi County Hospital (KCH, see
256 Supplementary Methods). The frequency of RSV occurrence among children <5 years of age
257 differed significantly (p-value <0.001) between the outpatient and inpatient setting during
258 the period of December 2015 to June 2017. The RSV viral load was similar between the
259 inpatient and outpatient settings (**Supplementary Figure 1(B)**). Across both inpatient and
260 outpatient care settings, the peak time for RSV case detections occurred from November to
261 May the following year (**Supplementary Figure 3**). The distribution of clinical symptoms
262 among patients attending the nine outpatient health facilities, and by age and gender are
263 shown in **Supplementary Table 2**.

264

265 **RSVB genome sequences from Kilifi.** WGS and data assembly was successful for
266 299/408 (73.3%) RSVB positive samples collected between December 2015 and June 2017
267 from the 9 selected health facilities in Kilifi county, Kenya. The remaining 26% (109/408)
268 sampling fraction were either of low viral load and we could not obtain more than 4 PCR
269 amplicons (81.6%) or were sequenced at insufficient depth or quality for genome assembly
270 (18.4%). 295/299 genomes were coding-complete (all the 11 RSV coding genes were
271 assembled) and were used in subsequent analysis. The median genome length was 15205
272 (range 11519 to 15257 nt). All the sequenced viruses belonged to the RSVB BA genotype,
273 characterized by the presence of 60-nt duplication in the C-terminal region of the G
274 glycoprotein gene. Genome coverage did not vary by rRT-PCR Ct value (**Figure 2(A)**).

275

276 **Genome-wide sequence diversity and evolution.** The sequence data showed linear
277 relationship between genetic change and time (root-to-tip correlation coefficient of 0.82,
278 **Figure 2(B)**); a temporal signal that supported the use of molecular clock models. Bayesian
279 analyses estimated an evolutionary rate of 1.063×10^{-3} (95% HPD: $9.2422 \times 10^{-4} - 1.2063 \times$
280 10^{-3}) nucleotide substitutions/site/year, and the genetic diversity traced back to a common
281 ancestor dated in 2014 (date in decimal format: 2014.89; 95% HPD: 2014.63 – 2015.14).
282 Demographic reconstruction (**Figure 2(C)**) showed seasonal periodicity in relative genetic
283 diversity that broadly mirrored RSVB incidence in the two epidemics. The peak genetic
284 diversity occurred in January followed by an inter-epidemic bottleneck that indicated
285 lineage or variant displacement events between epidemics.

286 The Kilifi genome sequences contained 838 consensus level non-gap single
287 nucleotide polymorphisms (SNPs), 554 (66%) of which were parsimony informative. 503/554
288 (91%) SNPs were located within coding regions, of which 332/503 (66%) were non-
289 synonymous and non-uniformly distributed across the RSV genome. Non-synonymous
290 changes peaked in density at the mucin-like domains of G gene; in the N-terminal of fusion
291 (F) gene; as well as in the N- and C-terminals of RNA-dependent RNA polymerase (L) gene
292 (**Figure 2(D)**).

293 Selection pressure analyses showed that G and SH glycoproteins had higher global
294 non-synonymous (dN)/synonymous (dS) substitution rate ratio estimates than all other
295 genes (**Table 2**). SLAC analyses identified three amino-acid sites (135, 217 and 285) in G
296 gene under significant positive selection ($P < 0.1$). In addition, we used MEME and FUBAR
297 methods to identify codons under pervasive and episodic (diversifying) positive selection.
298 MEME analyses detected 3 diversifying codons in the F gene, and 11 in the G gene ($P < 0.1$)

299 (Table 2). The FUBAR method identified 2 codon sites in F gene and 7 in G gene, under
300 episodic positive selection with significant support (posterior probability >0.9) (Table 2).

301

302 **Viral introductions and spread within Kilifi.** Genomic phylogenetic and
303 phylogeographic analyses afforded an in-depth look into the introduction and spread of
304 RSVB in the study population between December 2015 and June 2017. To determine viral
305 introductions, we analyzed the Kilifi RSVB genomes ($n=295$) in the context of other global
306 RSVB genomes ($n=500$). Sequences from Kilifi fell into eight major clades (numbered I to VIII,
307 **Figure 3(A)**) in which they clustered more closely among themselves and less within the
308 diversity of contemporaneous strains, indicating at least eight separate introductions into
309 the study population. The absence of external/global strains within the eight clades might
310 suggest local diversification of viruses, although we cannot exclude additional importation
311 from unsampled locations. Still, the low availability of RSVB sequence data from other
312 regions within Kenya or neighboring countries makes it difficult to estimate the precise
313 number of viral introductions.

314 The first four introductions (I-IV) comprised viruses confined to a single RSV
315 epidemic (**Figure 3(B)**). Clades I, II and IV comprised 25 viruses in total, from the 2015/16
316 epidemic, while clade III comprised 73 viruses uniquely from the 2016/17 epidemic. The
317 other four clades (V-VIII) comprised 76, 16, 36 and 69 viruses, respectively, a mix from both
318 RSV epidemics. All the eight introductions circulated within more than one location (table
319 inset, **Figure 3(A)**). Viruses in clade III in Figure 3(A) were closely related to RSVB strains
320 collected in January and August 2016 in UK and Australia, respectively. Clades I, IV, V and
321 VIII, shared common ancestors with sequences from Australia. Viruses in Clade VI were
322 closely related to a sequence collected in USA in 2015, while those in clade II were closely

323 related to a Dutch strain. There is a strong potential for bias in these results due to
324 heterogeneous sampling globally. We used TreeTime (42) for maximum likelihood dating of
325 the inferred clades, and the most likely dates of introduction were placed between July
326 2007 and September 2014 (**Supplementary Table 3**).

327 Analyzed separately, the time measured MCC tree of Kilifi genomes showed
328 phylogenetic clustering by epidemic, which depicted a temporally structured viral
329 population (**Figure 3(B)**). Our analysis of the geographical signal revealed high AI and PS
330 values (**Table 3**), suggesting relatively extensive viral migration dynamics among study
331 locations. The association between phylogenetic clustering and study location was
332 significant (p value < 0.001) in at least 8/9 study locations as revealed by the maximum
333 clade size values (**Table 3**). Differences in the observed and expected MC values (**Table 3**)
334 suggested that Mavueni exhibited the most spatial structure (difference of 8.7) and Mtondia
335 had the least (difference of 0.3).

336 We compared the model fit for symmetrical (reversible) and asymmetrical (non-
337 reversible) discrete diffusion models with BSSVS procedure. The asymmetrical CTMC model
338 gave a better fit (marginal likelihood, path sampling = -264.14 and stepping-stone sampling
339 = -316.3) compared to the reversible's marginal likelihood path sampling of -287.44 and
340 stepping-stone sampling of -351.9. The non-reversible model is a more realistic description
341 of the diffusion process that allows location exchange rates to vary according to
342 directionality (37). Mavueni may have played a central role in virus dissemination to other
343 study locations: it was the most probable location for most ancestral nodes with location
344 posterior support >0.9, (**Figure 4(A)**), and viral lineage movements from Mavueni were
345 statistically supported with Bayes Factor >1000 (**Table 4**). In addition, there was a relatively
346 high genomic diversity in Mavueni (red-colored taxa in **Supplementary Figure 4**). The

347 inferred viral movements between the nine study locations are visualized as geographic
348 links in **Figure 4(B)**. There were very few statistically supported viral movements between
349 Sokoke and other study locations.

350

351 **Phylogenomic clusters.** The RSVB viruses in Kilifi were further summarized as
352 phylogenomic clusters based on genetic similarity of consensus genomes. For this, we
353 examined the distribution of maximum likelihood pairwise evolutionary distances among
354 the 295 samples (43365 pairs). The genetic distance distribution was multi-modal (**Figure 5**
355 **(A)**), and we defined a phylogenomic cluster as a group of viruses (≥ 2) within an epidemic
356 that contained evolutionary pairwise distances of <0.0024 nucleotide substitutions per site
357 (red dashed line **Figure 5 (A)**) or 37 pairwise nucleotide differences. The distance threshold
358 of <0.0024 nucleotide substitutions per site demarcated the 17th percentile of the entire
359 distance distribution and was additionally supported by clear separation of viruses into
360 closely related phylogenetic clusters with >0.9 posterior probability support.

361 Overall, at least 20 phylogenomic clusters were identified, all of which were
362 supported by strong posterior probability support (>0.9); nine in the 2015/16 epidemic
363 (**Figure 6(A)**) and eleven (**Figure 6(B)**) in the 2016/17 epidemic. The lowest within-cluster
364 genetic diversity was found in cluster 12 in the 2016/17 epidemic and the highest in cluster
365 four in the 2015/16 epidemic (**Figure 5(B)**). Genetically, cluster five in the 2015/16 epidemic
366 fell basal to cluster 17 in the 2016/17 epidemic and both made the 6th RSVB introduction
367 (clade VI in **Figure 3(A)**), a phylogenetic pattern compatible with in-situ virus evolution and
368 persistence. The clusters ranged in size, from 2 to 72 viruses. Overall, there was abundant
369 viral diversity in Pingilikani in both epidemics (13 clusters), while Matsangoni and Ngerenya

370 had the lesser diversity (six clusters). In particular, only one cluster or viral lineage was
371 circulating in Matsangoni in the 2015/16 epidemic.

372 Sample collection dates were used to estimate the duration of each phylogenomic
373 cluster. Clusters 20 (172 days), 10 (173 days) and six (211 days) circulated over a wide
374 temporal scale up to several months. Geographically, clusters one and 20 were the largest in
375 size and most pervasive circulating widely in all the nine study locations (**Figure 6**). Two
376 viruses in the 2016/17 epidemic (marked with '*' in **Figure 6**) were singletons and could not
377 be assigned to any phylogenomic cluster, likely an artifact of limited sampling. The number
378 of clusters identified in each study location was weakly associated with the number of
379 genomes from that location (Spearman correlation $\rho = 0.2$, $P = 0.6$)

380

381 **Amino acid diversity in F and G genes.** Given their role in the initial phases of
382 infection and as the major antigens for eliciting neutralizing antibody responses, we
383 examined amino-acid variation in the two major glycoproteins on the surface of the RSV
384 virion, the attachment (G) glycoprotein and the fusion (F) glycoprotein. The F protein is also
385 the major target for antiviral immuno-prophylaxis. In all the sequenced viruses, there were
386 37 and 93 amino-acid (codon) sites containing non-synonymous changes in F and G gene,
387 respectively, and the frequency of the observed non-synonymous variants was up to 24.7%
388 and 90.4% in F and G gene, respectively. While most of these variable sites contained a
389 single non-synonymous change, one codon in F contained two changes (K419N/R), and
390 three codons in G contained multiple non-synonymous changes (T107A/D, P289L/S and
391 E290A/G/V). The K419N and K419R substitutions were seen in a subset (n=6) of sequences
392 in cluster 13 and in cluster 17, respectively. The T107D substitution was observed in a
393 subset (n=7) of sequences in cluster 13, while the T107A was seen in all other clusters

394 except for seven, eight and nine. The P289L and P289S substitutions were observed in
395 cluster seven and eight, and six, respectively. The E290A, E290G and E290V substitutions
396 were observed in clusters 13, 6 and 4 and 14, respectively. These amino-acid variants were
397 not limited to a study location and were relative to the earliest dated sequence in the
398 dataset (14 December 2015).

399 In F gene, viruses in clusters 10 and 11 (**Figure 6(B)**) in the 2016/17 epidemic
400 contained two unique amino-acid changes (N99T and T129I). Additionally, three non-
401 synonymous substitutions (V103A, L172Q and S173L) located in the F protein antigenic
402 site V were unique to clusters seven, eight and nine in the 2015/16 epidemic. In G gene,
403 90% of the genomes contained T107A non-synonymous substitution, 64% had P217L
404 substitution, and 45% contained a H285Y substitution, compared to the earliest sequence.
405 A number of amino-acid changes occurred at the node between cluster nine and other
406 viruses (marked A in **Figure 3(A)**) and these included N213D, A227I, P255S, and N303E in G
407 gene; and T19A and P125L in F gene. A non-synonymous change from CAA(Q) resulted in
408 early Stop codon (TAA) at amino acid position 311 in 135/295 (45%) sequenced viruses.

409 Several amino-acid sites appeared to undergo reversal mutation, for instance in G
410 gene, T/S75S/T occurred within the 2015/16 epidemic, while T/I252I/T, L/P270P/L and
411 G/S135S/G occurred in both epidemics. The I/L542L/I reversal mutation in F gene was seen
412 in the 2015/16 epidemic. The T/S75S/T (G gene) and I/L542L/I (F gene) reversal
413 substitutions were both observed in phylogenomic clusters seven and eight that comprised
414 clade II (**Figure 3(A)**). The L/P270P/L reversal mutation occurred in clusters 5, 7 and 17.
415 Reversible evolution may contribute to the escape from the human population immune
416 response, thereby facilitating viral transmission (43).

417

418 **Emerging genetically distinct variant in 2016/17 epidemic.** Several non-synonymous
419 changes were unique to phylogenomic cluster 20 in the 2016/17 epidemic (**Figure 6(B)**),
420 including K68Q in F gene; Y90H, L91F, T225N, T273I, and A301T in G gene. These six non-
421 synonymous changes were co-occurring in all the sequences in cluster 20 and were detected
422 at an intermediate frequency (41% of samples in 2016/17) in the study population.
423 Importantly, the amino-acid change K68Q is at the antigenic site \emptyset , the binding site of
424 monoclonal antibody (mAb) MEDI8897. A variant with the mutation K68N was previously
425 detected in 2% of sampled viruses in the US (44). Other amino-acid changes distinctive to
426 the cluster 20 are listed in **Supplementary Table 4**. This phylogenomic cluster, which was
427 also unequivocally an independent introduction (clade III, **Figure 3(A)**), was the most
428 prevalent and its detection through genomic analysis rules out any sample contamination
429 and sequencing errors.

430 We went further and checked whether this emerging variant was also present in the
431 wider Kilifi county based on G gene sequence data of samples collected from pediatric ward
432 admissions (<59 months of age) at the Kilifi County Hospital (KCH, formerly Kilifi District
433 Hospital) during the 2016/17 RSV epidemic. KCH provides primary care and inpatient
434 referral services to a larger catchment area. We found that 32% of RSVB positive hospital
435 samples in the 2016/17 epidemic clustered with the emerging variant and likewise
436 contained the non-synonymous changes Y90H, L91F, T225N, T273I, and A301T in G gene,
437 which further increased the overall prevalence of clade III in the 2016/17 epidemic.

438 We also investigated whether the emerging variant persisted in Kilifi after the
439 2016/17 RSV epidemic, based on G gene sequences of RSVB positive samples collected from
440 pediatric admissions at KCH during the 2018/19 epidemic. To note, there were no RSVB
441 cases in the 2017/18 epidemic in the KHDSS study (outpatient surveillance in the nine health

442 facilities) (**Supplementary Figure 5**) and only 2 RSVB cases were detected in the in-patient
443 pediatric admissions at KCH (<0.01% of RSV positive samples). We found that viruses in
444 2018/19 epidemic clustered closely with clade III or cluster 20 and the branch leading to the
445 2018/19 RSVB epidemic shared a common ancestor with cluster 20 (**Supplementary Figure**
446 **6**). The 2018/19 viruses likely descended from the same introduction as clade III, suggesting
447 that the mutations defining clade III were epidemiologically important. It is possible that
448 RSVB underwent a significant population (or genetic) bottleneck after the 2016/17 epidemic
449 and older viruses were eliminated given that no RSVB viruses or persistent variants were
450 transmissible in the 2017/18 epidemic.

451

452 **Amino-acid diversity in other regions of the genome.** Additional variable amino acid
453 sites observed in the other genomic regions are listed in (**Supplementary Table 4**).
454 Prominent amino acid changes occurring in more than one phylogenomic cluster are
455 underlined in **Supplementary Table 4**. The emerging variant (cluster 20) discussed above
456 had two other unique amino acid differences in NS2 (non-structural protein 2) and L
457 (polymerase) genes. Even in other genomic regions, phylogenomic clusters 5 (2015/16
458 epidemic) and 17 (2016/17 epidemic), which are nested within the same virus introduction
459 (clade VI, **Figure 3(A)**), shared amino-acid differences. Clusters 10 to 14 had the same
460 amino-acid differences in F and L genes.

461

462 **Discussion**

463 Here we report whole genome molecular epidemiology and phylodynamics of
464 respiratory syncytial virus group B providing a detailed view of the introduction and spread
465 the virus in Kilifi county, coastal Kenya. The results were obtained from genomic analyses of

466 295 samples originating from representative sampling across the KHDSS area, over two
467 consecutive RSV epidemics. Phylogenetic analyses revealed multiple virus introductions,
468 each introduction commonly circulating in all the study locations, suggesting substantial
469 spatial spread and transmission between locations in a relatively short time. Although RSV
470 surveillance has improved in many regions across the world, publicly available RSVB
471 genomic data from recent years is quite insufficient and may have limited our inference of
472 spatial and temporal placement of virus introductions in Kilifi.

473 Analysis of RSVB transmission dynamics in Kilifi suggested extensive viral migrations
474 among the study locations associated with the nine health facilities as well as strong spatial
475 substructures within each location. The substructures might represent predominant local
476 transmission and diversification processes. Using the BSSVS approach and an asymmetrical
477 diffusion model to reconstruct RSVB dispersal, we observed significant and strong support
478 for epidemiological links between one central location (Mavueni), located at the
479 intersection of the main roads through the KHDSS and other study locations (BF > 1000,
480 Table 4). However, we note that Mavueni had a larger proportion of RSVB genomes
481 compared to other locations, which might bias the ancestral location estimates. Improved
482 road infrastructure and transportation within KHDSS has facilitated mobility thus
483 increasingly connecting the local population and expanding virus transmission networks.

484 Temporal sequence divergence and accumulation of nucleotide substitutions was
485 detectable over the sampling timeframe and varied by epidemic intervals as shown by two
486 separate groups of tips in the linear regression plot (Figure 2B). Time-scaled phylogenies
487 also exhibited a marked epidemic behavior indicative of chronological generation of new
488 variants. The two epidemics were characterized by multiple clades for viruses sampled
489 within the same epidemic, indicating continued transmission generated and sustained by

490 increasing spatial connectivity in the wider Kilifi county. Changes in the relative genetic
491 diversity coincided with RSVB case detection and captured fine temporal resolution of
492 changes in the viral population size, which also implied sufficient sampling density (14).
493 While estimating nucleotide substitution rate is useful in revealing the dynamics and
494 processes of viral evolution (45), we could not directly compare our estimate of the mean
495 nucleotide substitution rate with previous studies due to varying sampling timescales,
496 different molecular clock (fixed vs. relaxed) and coalescent models, as well as epidemiologic
497 variations.

498 Phylogenetic clusters have been used to investigate epidemiologically significant HIV
499 hotspots and characterize groups burdened by a high rate of HIV transmission (46). We
500 inferred 20 phylogenomic clusters of closely related viruses based on genetic relatedness,
501 which we put forward as potential transmission units. Stochastic difference in
502 transmissibility or circulation, or infection rates could explain the differences in prevalence
503 of the different phylogenomic clusters.

504 Using SLAC, MEME, and FUBAR methods, there was limited detectable diversifying
505 selection for the amino acid substitutions that characterized or defined the different
506 phylogenomic clusters. Only three cluster-defining codons in G (144, 294 and 303) and three
507 in F (125, 172 and 173) glycoproteins were subject to episodic positive selection. It is likely
508 that the observed codon replacements follow non-selective epidemiological processes and
509 these substitutions are compensatory mutations to retain function, or hitchhikers carried
510 along by chance (47), or that immune driven positive selection could not be identified by the
511 three methods used in this study. Nonetheless, the implications of these mutations on
512 protein function, viral evolution and fitness are uncertain and warrant further functional
513 investigation.

514 A previous study showed that viruses carrying the K68N substitution in the US,
515 affected the binding of MEDI8897 (48). MEDI8897 is an RSV pre-F-specific human mAb with
516 an extended serum half-life, under clinical evaluation as a passive immunization of all
517 infants entering their first RSV season (49). It is probable the K68Q substitution identified in
518 Kilifi promoted evasion of pre-existing immune responses. The unequivocal support for the
519 monophyly of the cluster 20 as an introduction into Kilifi (clade III, **Figure 3(A)**), further
520 supports that this phylogenetic clade was a single (epidemiologically successful)
521 introduction event. Surprisingly, the K68Q F gene amino-acid substitution observed in
522 phylogenomic cluster 20 was not detected as under any selection pressure even though this
523 residue is located at structurally determined mAb binding epitopes (50). Perhaps an
524 explanation for this is the low rate of nonsynonymous evolution (conversely, high sequence
525 conservation) at codon 68 in our dataset. In any case conventional approaches for
526 measuring selection pressure consistently detect positive selection only at codon sites with
527 high rates of nonsynonymous evolution (51).

528 A high sequence conservation was reported in the MEDI8897 binding site among
529 naturally occurring RSV isolates collected from 1965 to 2014 (50). Our study provides a
530 novel sequence polymorphism (K68Q) within the MEDI8897 binding site with a frequency of
531 nearly 50% in the 2016/17 RSV epidemic in our study population. Functional
532 characterization is required to determine MEDI8897 neutralization and/or binding to viruses
533 containing the K68Q mutation. Additionally, the viruses with the K68Q mutation in 2016/17
534 epidemic possessed five distinctive amino-acid mutations in G gene, including Y90H and
535 L91F in two consecutive codons. We cannot exclude the possibility that these are relevant
536 antigenic epitopes. Our study underscores the need for continued genomic surveillance of

537 contemporary clinical strains particularly at F and G protein antigenic sites as this has
538 implications on RSV therapeutic and vaccine development.

539 RSVB viruses containing A103V/L172Q/S173L amino acid-changes in the F protein
540 were also detected during the 2015/16 RSV season in USA and estimated to have likely
541 emerged around 2014 (44) and in China (52), suggesting global circulation of this variant.
542 However, unlike in the US, none of the samples from 2016/17 in Kilifi had these three
543 substitutions probably due to removal by purifying selection.

544 In conclusion, we present the utility of genomic analyses to investigate virus
545 transmission and genetic diversity including detection of a novel antigenically distinct
546 variant. Further studies are required to determine whether the K68Q mutation is adaptive
547 and/or a result of escape from antibody-mediated selection and constitutes a naturally
548 acquired antiviral resistance-associated mutation that potentially disrupts neutralizing
549 antibody recognition and binding. An important future surveillance effort for us is to assess
550 if the K68Q mutation has become more prevalent and gradually fixed since the 2016/17
551 epidemic. Additional sequencing of RSVB from other regions in Kenya and neighboring
552 countries is essential to refine evolutionary dynamics and draw better conclusions about
553 geographic origins of viral introductions in the study population in Kilifi. The present study
554 makes publicly available a large number of newly acquired coding-complete RSVB genomes
555 useful for further molecular evolution studies.

556

557 **Data availability**

558 The replication data and analysis scripts for this manuscript are available from the Harvard
559 Dataverse: [???](#). Some of the clinical dataset contains potentially identifying information on

560 participants and is stored under restricted access. Requests for access to the restricted
561 dataset should be made to the Data Governance Committee (dgc@kemri-wellcome.org).

562

563 **Competing interests**

564 The authors declare no competing interests.

565

566 **Acknowledgements**

567 We thank all the study participants for their contribution of samples and data. We also
568 thank the Dispensary / Health Centre management committees for allowing us to conduct
569 the study within their health facilities. We are grateful to the field study team for participant
570 recruitment and the laboratory staff of the KEMRI-Wellcome Trust Research Programme /
571 Virus Epidemiology and Control research group. We would also like to thank D. Collins
572 Owuor for his assistance with the BaTS software and Mark Otiende for his assistance with
573 demographics data. This paper is published with the permission of the Director of KEMRI.

574

575 **Funding**

576 This work was supported by the Wellcome Trust [grant 102975, 203077]. CNA is supported
577 by the Initiative to Develop African Research Leaders (IDeAL) through the DELTAS Africa
578 Initiative [DEL-15-003]. The DELTAS Africa Initiative is an independent funding scheme of the
579 African Academy of Sciences (AAS)'s Alliance for Accelerating Excellence in Science in Africa
580 (AESA) and supported by the New Partnership for Africa's Development Planning and
581 Coordinating Agency (NEPAD Agency) with funding from the Wellcome Trust
582 [107769/Z/10/Z] and the UK government. The views expressed in this publication are those

583 of the authors and not necessarily those of AAS, NEPAD Agency, Wellcome Trust or the UK
584 government.

585
586 **Figure legends**

587 **Figure 1(A)** A map showing the geographical area covered in the Kilifi Health Demographic
588 Surveillance System (KHDSS), expanded from a map of Kenya. The nine participating public
589 health facilities are indicated in the map. The dark lines within the polygons indicate the
590 road structure within KHDSS. The maps were rendered using QGIS 2.18.17

591 (<https://www.qgis.org/>) **(B)** Monthly RSVB occurrence by study location: temporal and
592 spatial distribution of RSVB positive cases (left Y axis) and number of clinical samples tested
593 (right Y axis) from each participating health facilities. Abbreviations: CHA = Chasimba, JAR =
594 Jaribuni, JUN = Junju, MAV = Mavueni, MAT = Matsangoni, PIN = Pingilikani, NGE =
595 Ngerenya, SOK = Sokoke, MTO = Mtondia

596
597 **Figure 2(A)** Genome coverage for each virus isolate versus the viral load (rRT-PCR Cycle
598 threshold (Ct) value). **(B)** Root-to-tip divergence as a function of sampling time colored by
599 study location. **(C)** Relative genetic diversity through time estimated using the Gaussian
600 Markov Random Field (GMRF) Skygrid model. Solid line represents mean relative genetic
601 diversity while the corresponding dashed lines indicate the 95% HPD intervals. **(D)** Relative
602 frequencies of potential non-synonymous changes across codon-aligned RSV genome
603 sequences. The frequencies for each codon position are calculated as the number of non-
604 synonymous nucleotide substitutions for all pairwise comparisons in a sequence alignment,
605 while excluding ambiguous bases or insertions. Abbreviations: CT = cytoplasmic, TM =
606 transmembrane, CCD = central conserved domain; SP = signal peptide; RdRp = RNA

607 dependent RNA polymerase, Cap = capping, and MT = methyltransferase, CD = connector
608 domain, CTD = C-terminal domain.

609

610 **Figure 3(A)** Phylogenetic placement of RSVB genomes from Kilifi relative to sequences from
611 other geographical regions. The clade or introduction assignments are indicated as I to VIII
612 on the trees and the circles at the branch tips are colored by location or country of origin.
613 Only the contemporaneous (global) sequences that are closely related to or share common
614 ancestors with RSVB genomes from Kilifi, are colored. The black circles represent other
615 global sequences used in the phylogenetic analysis. Geographic distribution of the RSVB
616 introductions in Kilifi is shown in the inset table. Abbreviations: CHA = Chasimba, JAR =
617 Jaribuni, JUN = Junju, MAV = Mavueni, MAT = Matsangoni, PIN = Pingilikani, NGE =
618 Ngerenya, SOK = Sokoke, MTO = Mtondia. **(B)** Maximum likelihood phylogeny (1000
619 bootstrap resampling) showing the temporal distribution of the RSVB introductions in Kilifi.
620 The individual sequences from Kilifi are colored by respective RSV epidemic. The black
621 circles represent other global sequences used in the phylogenetic analysis. Bootstrap
622 support values are shown for the most basal nodes of the inferred introductions.

623

624 **Figure 4(A)** Temporal-scaled phylogeographic maximum clade credibility tree of Kilifi RSVB
625 genomes. Branch colors indicate the modal location (most probable reconstructed location
626 at each node) inferred under the asymmetrical discrete phylogeographic model. Location-
627 state posterior probabilities are shown next to relevant nodes along with the posterior
628 support. **(B)** Spatial diffusion pathways of RSVB in Kilifi. Well supported (Bayes Factor >6)
629 discrete diffusion asymmetric rates of viral movement between the study locations are

630 indicated. The type and thickness of the arrows represents the relative strength of the
631 diffusion rate.

632

633 **Figure 5(A)** Histogram of the whole genome sequences patristic distance frequency
634 distribution. The vertical red dashed line corresponds to the 17th percentile distance
635 threshold (0.0024 expected nucleotide substitutions per site) for which phylogenomic
636 clusters were identified. The distances were measured in units of nucleotide substitutions
637 per site and extracted from a maximum likelihood phylogeny (1000 bootstrap resampling).
638 **(B)** Maximum pairwise patristic distances for sequences belonging to the different
639 phylogenomic clusters.

640

641 **Figure 6** The assigned phylogenomic clusters are shown on maximum likelihood phylogenies
642 (1000 bootstrap resampling): 9 clusters in the 2015/16 epidemic **(A)** and 11 clusters in the
643 2016/17 epidemic **(B)**. Geographical composition and number of viruses/sequences per
644 phylogenomic cluster are shown (inset tables). The '*' in (B) indicates sequences that were
645 not assigned into a phylogenomic cluster. All sequences had epidemiological and
646 geographical information. Amino-acid changes in F and G glycoproteins unique to the
647 emergent variant (cluster 21) in the 2016/17 are labeled on the major branch leading to the
648 cluster. Scale bar represents the number of substitutions per site.

649

650 **Supplementary Figure 1 (A)** Violin plots show the distribution (median, IQR) of RSVB qRT-
651 PCR Cycle threshold (C_t) values (≤ 40) for outpatient participants. **(B)** Comparison of C_t
652 distribution between inpatient (KCH) and the outpatient facilities for children < 5 years for

653 the period of December 2015 to June 2017. The threshold used for determining positive and
654 negative samples is shown by dashed line ($C_t=35.0$).

655

656 **Supplementary Figure 2** Frequency (primary Y axis) and proportion (secondary Y axis) of
657 RSV across different age categories in both genders (dashed and dotted lines) from
658 participants presenting to the nine outpatient health facilities over the period December
659 2015 to June 2017.

660

661 **Supplementary Figure 3** Distribution by month of RSVB positive samples collected from the
662 inpatient (severe pneumonia admissions to KCH) and outpatient facilities over the period of
663 December 2015 to June 2017. Secondary Y axis records sample tested (dotted line) and
664 proportion (dashed line) that is RSVB positive per month.

665

666 **Supplementary Figure 4** Temporal-scaled phylogeographic maximum clade credibility tree
667 showing sequences of samples collected in Mavueni location (red taxa) in Kilifi County,
668 coastal Kenya.

669

670 **Supplementary Figure 5** The number of RSVB cases through time (2015 to 2019) for severe
671 pneumonia pediatric admissions to Kilifi County Hospital. RSV epidemics (usually from
672 October to July of the following year) are indicated.

673

674 **Supplementary Figure 6** Maximum likelihood phylogeny of the attachment G glycoprotein
675 gene. The orange colored branches represent the recently emergent variant (phylogenomic

676 cluster 20). The purple colored branches represent the 2018/19 RSVB viruses collected from

677 severe pneumonia pediatric admissions to Kilifi County Hospital.

678

679

680 **Tables**

681 **Table 1** Characteristics of RSV positive and negative study participants by age, gender, and
682 the 9 outpatient health facilities in Kilifi County, coastal Kenya.

Characteristic	NPS virus positive(n)	%	NPS virus negative (n)	%	Total(n)	P value
	(n=503)	6.19	7624	93.8	8127	
Age in years						
Mean			149.4		144.0	
Median (IQR*)	20 (8-43)		54(18-194)		49 (17-188)	
Sex						
Male	231	6.7	3195	93.3	3426	0.078
Female	272	5.8	4428	94.2	4700	
Age Category						
0–5 mo	85	11.3	665	88.7	750	<0.001
6–11 mo	84	11.0	681	89.0	765	
12–23 mo	102	9.1	1021	90.9	1123	
24–35 mo	76	9.6	712	90.4	788	
3–4 y	64	6.7	898	93.4	962	
5–9 y	37	3.4	1046	96.6	1083	
10–19 y	21	2.0	1024	98.0	1045	
20–49 y	23	2.1	1098	98.0	1121	
50–100 y	11	2.3	479	97.8	490	
Health Facility						
Matsangoni	58	6.1	895	93.9	953	0.003
Ngerenya	58	6.5	840	93.5	898	
Sokoce	50	5.7	826	94.3	876	
Mtondia	66	6.9	898	93.2	964	
Mavueni	79	8.6	838	91.4	917	
Jaribuni	43	5.1	803	94.9	846	
Chasimba	67	7.5	826	92.5	893	
Junju	40	4.3	889	95.7	929	
Pingilikani	42	4.9	809	95.1	851	

683

684 **Table 2** The predicted nature of selection pressures acting on each genomic region: 1st
685 column shows the computed mean dN/dS rate ratio using SLAC and the 2nd column shows
686 amino-acid sites in F and G gene under episodic selection as identified by MEME analyses.
687 Sites also detected using the FUBAR method, in addition to MEME, are underlined.

Non-synonymous (dN)/synonymous (dS) substitution rate ratio per site		Sites subject to episodic positive/diversifying selection
NS1	0.12	<i>G gene</i>
NS2	0.236	<u>135</u> , 144, <u>154</u> , 172, 208, <u>217</u> ,
N	0.0832	<u>285</u> , 291, <u>294</u> , 298, <u>303</u>
M	0.0525	<u>252</u> *
P	0.0642	
F	0.179	<i>F gene</i>
G	0.487	<u>125</u> , <u>172</u> , 173
SH	0.426	
M2-1	0.264	
M2-2	0.267	
L	0.122	

688 * This site was detected by the FUBAR method only.

689

690

691 **Table 3** Results of Bayesian analysis of phylogeographic structure of RSVB viruses in Kilifi,
 692 coastal Kenya, 2015-2017. *P* values correspond to the proportion of trees from the expected
 693 (null) distribution equal to, or more extreme than, the median posterior of the statistic.
 694 Abbreviations: CHA = Chasimba, JAR = Jaribuni, JUN = Junju, MAV = Mavueni, MAT =
 695 Matsangoni, PIN = Pingilikani, NGE = Ngerenya, SOK = Sokoke, MTO = Mtondia.

Location	Association Index (AI) (95% CI) †			Parsimony Score (PS) (95% CI) †			Maximum Clade size (95% CI) ‡			Difference §
	Observed	Expected	<i>P</i> value	Observed	Expected	<i>P</i> value	Observed	Expected	<i>P</i> value	
ALL	14.8 (13.8-15.8)	31 (29.5-32.3)	0.0	132.6 (129-136)	208.4 (202-214.2)	0.0	-	-	-	-
CHA	-	-	-	-	-	-	4.6 (4-6)	1.84 (1.1-2.7)	10E-4	2.62
JAR	-	-	-	-	-	-	5.27 (5-6)	1.56 (1-2.08)	10E-4	3.71
JUN	-	-	-	-	-	-	6 (6-6)	1.55 (1-2.03)	10E-4	4.45
MAT	-	-	-	-	-	-	6.32 (4-10)	1.71 (1-2.2)	10E-4	4.61
MAV	-	-	-	-	-	-	11 (11-11)	2.3 (1.76-3)	10E-4	8.7
MTO	-	-	-	-	-	-	2 (2-2)	1.7 (1-2.3)	0.21	0.3
NGE	-	-	-	-	-	-	3.65 (2-4)	1.64 (1-2.2)	10E-4	2.01
PIN	-	-	-	-	-	-	3.1 (3-4)	1.73 (1-2.2)	0.0084	1.37
SOK	-	-	-	-	-	-	4 (3-5)	1.8 (1.1-2.3)	10E-4	2.2

696 † AI and PS metrics were determined for all locations combined.

697 ‡ Maximum clade size was determined for each specific location.

698 § Difference between observed and expected (null) clade size.

699

700 **Table 4** Statistically supported state transitions indicating viral migration events. Bayes
701 factor >6 was considered significant.

Transition between		Distance (Km) ^Φ	Bayes factor	Mean indicator ^δ	Mean rate
Chasimba	Mtondia	21.5	6.2	0.46	1
Jaribuni	Junju	17.5	6.7	0.48	1
Chasimba	Ngerenya	25.5	8.3	0.53	0.8
Sokoce	Matsangoni	20.6	9.8	0.57	0.9
Mtondia	Junju	24	12.2	0.63	0.76
Chasimba	Pingilikani	6.2	13.7	0.65	0.9
Pingilikani	Mtondia	21.6	18.2	0.71	1
Junju	Jaribuni	17.5	20.8	0.74	1
Mavueni	Mtondia	12.5	28.5	0.80	1
Matsangoni	Ngerenya	12.3	35.0	0.83	0.75
Mavueni	Matsangoni	29.3	179.9	0.96	0.87
Mavueni	Junju	11.5	434.6	0.98	0.87
Junju	Pingilikani	5.4	1166.6	0.99	1.4
Ngerenya	Mtondia	4.6	1624.5	1.00	1.3
Mavueni	Ngerenya	17	2062.5	1.00	1.2
Sokoce	Mtondia	10.9	12678.4	1.00	1.7
Mavueni	Jaribuni	12.8	24571.2	1.00	1.09
Mavueni	Pingilikani	9.8	30243.2	1.00	1.5
Mavueni	Chasimba	12.1	196620.8	1.00	1.75
Mavueni	Sokoce	16.5	393248.8	1.00	1.78

702 ^Φ Great circle distance estimates between centroid longitude and latitude of the respective locations.

703 ^δ The posterior probability of observing non-zero migration rates in the Bayesian sampled trees.

704

705

706 **Supplementary Table 1** Information (accession number, country of sampling and collection

707 date) on RSVB genome sequences of >14,000 nucleotide length retrieved from NCBI

708 GenBank.

709

710

711 **Supplementary Table 2** Frequency of distribution of symptoms among study participants by

712 9 outpatient health facilities, age, and gender, between December 2015 and June 2017.

Characteristic	Fever, %	Chest indrawing, %	Crackles, %	Wheeze, %	Cough, %	Nasal discharge, %	Nasal flare, %	Difficulty breathing, %	Total participants, n
Health Facility									
Matsangoni	82.8	0.0	3.5	6.9	98.3	72.4	3.45	10.3	953
Ngerenya	84.5	1.7	1.7	0.0	98.3	65.5	0.0	10.3	898
Soko	84.0	2.0	2.0	6.0	100	80.0	4.0	18.0	876
Mtondia	81.8	6.1	3.0	1.5	97.0	78.8	4.6	7.6	964
Mavueni	74.7	3.8	2.5	1.3	97.5	77.2	5.1	15.2	917
Jaribuni	79.1	2.3	4.7	9.3	97.7	83.7	0.0	4.7	846
Chasimba	65.7	3.0	1.5	6.0	97.0	71.6	3.0	9.0	893
Junju	80.0	0.0	5.0	7.5	95.0	90.0	2.5	20.0	929
Pingilikani	85.7	2.4	2.4	9.5	95.2	76.2	4.8	40.5	851
P value	0.156	0.588	0.974	0.142	0.900	0.192	0.741	0.001	
Age Category									
0–5 mths	77.7	3.5	2.4	5.9	98.8	69.4	10.6	23.5	750
6–11 mths	83.3	8.3	4.8	6.0	100	82.1	3.6	22.6	765
12–23 mths	84.3	2.0	3.9	5.9	97.1	73.5	1.0	13.7	1123
24–35 mths	82.9	1.3	2.6	4.0	96.1	81.6	2.6	9.2	788
3–4 yrs	79.7	0.0	1.6	6.3	98.4	76.6	1.6	9.4	962
5–9 yrs	83.8	0.0	2.7	2.7	91.9	83.8	0.0	8.1	1083
10–19 yrs	66.7	0.0	0.0	0.0	100	66.7	0.0	4.8	1045
20–49 yrs	52.2	0.0	0.0	0.0	91.3	87.0	0.0	4.4	1121
50–100 yrs	45.5	0.0	0.0	0.0	100	54.6	0.0	0.0	490
P value	0.003	0.041	0.883	0.830	0.129	0.150	0.009	0.010	
Sex									
Male	80.5	3.5	3.0	5.6	98.7	73.2	2.2	16.5	3426
Female	77.9	1.8	2.6	4.0	96.3	79.4	4.0	12.1	4700
P value	0.478	0.252	0.756	0.406	0.094	0.099	0.231	0.166	

713

714 **Supplementary Table 3** Rate of nucleotide substitution and the estimated date of
715 introduction or emergence of each clade.

Clade	Estimated date	Rate of evolution (10^{-3} subs/site/y)		
		Mean	Lower 95% HPD	Upper 95% HPD
I	September 2008	1.2	0.27	2.3
II	July 2007	1.18	0.33	2.32
III	September 2014	1.06	0.35	1.95
IV	August 2012	0.43	0.26	0.54
V	March 2013	1.22	0.87	2.02
VI	August 2012	0.67	0.58	1.1
VII	May 2013	1.64	1.03	2.17
VIII	March 2014	0.92	0.41	1.87

716

717

718 **Supplementary Table 4** Amino acid substitutions identified in RSV genomic regions, which
719 were distinctive to the various phylogenomic clusters.

Cluster (2015/16)	Amino acid substitutions in genomic regions	Cluster (2016/17)	Amino acid substitutions in genomic regions
1	<i>P gene</i> (Ile ₆₀); <i>L gene</i> (Ile ₁₀₅ , Asn ₁₈₃)	10	<i>F gene</i> (Val ₅ , Thr ₉₉ , Thr ₁₂₉); <i>G gene</i> (Ser ₁₀₁ , Phe ₂₆₇); <i>NS2 gene</i> (Asn ₅₃); <i>L gene</i> (Asp ₁₇₆)
2	<i>N gene</i> (Val ₉₇); <i>G gene</i> (His ₁₄₄)	11	<i>F gene</i> (Val ₅ , Thr ₉₉ , Thr ₁₂₉); <i>L gene</i> (Asp ₁₇₆ , Leu ₁₇₀₈); <i>G gene</i> (Phe ₂₆₇)
3	<i>N gene</i> (Val ₉₇)	12	<i>F gene</i> (Ser ₁₀₄ , Val ₅); <i>G gene</i> (Phe ₂₆₇); <i>NS2 gene</i> (Arg ₁₀₁); <i>L gene</i> (Asp ₁₇₆)
4	<i>G gene</i> (Val ₂₉₀)	13	<i>F gene</i> (Val ₅); <i>L gene</i> (Asp ₁₇₆); <i>G gene</i> (Thr ₂₆₈ , Ala ₂₉₀ , Leu ₃₀₄)
5	<i>NS2 gene</i> (Thr ₅); <i>L gene</i> (Ile ₁₁₆₆ , Arg ₂₀₆₆); <i>G gene</i> (Val ₂₆₉ , Pro ₂₇₀)	14	<i>F gene</i> (Val ₅); <i>L gene</i> (Asp ₁₇₆ , Phe ₁₉₈₀); <i>G gene</i> (Thr ₂₆₈ , Val ₂₉₀)
6	<i>F gene</i> (Thr ₁₆); <i>L gene</i> (Ala ₁₇₄₄ , Gly ₁₇₈₇)	15	<i>N gene</i> (Val ₉₇); <i>G gene</i> (Thr ₂₀₅ , Leu ₂₁₄ , Ala ₂₄₉ , Leu ₂₉₃ , Pro ₃₀₅)
7	<i>F gene</i> (Leu ₁₇₂ , Ser ₁₇₃ , Leu ₅₄₂); <i>G gene</i> (Ser ₇₅ , Asn ₁₀₃ , Thr ₁₀₇ , Arg ₁₃₆ , Gly ₂₆₀ , Pro ₂₇₀ , Ile ₂₇₉ , Leu ₂₈₉ , Lys ₃₁₂); <i>M2-1 gene</i> (Val ₁₈₁);	16	<i>F gene</i> (Cys ₂₅); <i>G gene</i> (His ₁₄₄); <i>L gene</i> (Glu ₃₀₄ , Ile ₁₉₅₆)

	<p><i>M2-2 gene</i> (Asn₂₆, Lys₄₉);</p> <p><i>N gene</i> (His₂₁₆);</p> <p><i>NS1 gene</i> (Asp₆₉, Ile₇₂);</p> <p><i>NS2 gene</i> (Thr₂₆);</p> <p><i>SH gene</i> (Thr₄₉);</p> <p><i>L gene</i> (Leu₅₆, Val₇₁₅, Ala₁₇₁₂, Asn₁₇₃₆, Ile₂₀₁₉, Ile₂₀₆₉)</p>		
8	<p><i>G gene</i> (Ser₇₅, Ile₇₆, Asn₁₀₃, Thr₁₀₇, Leu₁₁₀, Arg₁₃₆, Ser₁₇₆, Leu₁₉₁, Ser₂₃₅, Gly₂₆₀, Ile₂₇₉, Leu₂₈₉, Lys₃₁₂);</p> <p><i>F gene</i> (Ala₁₀₃, Thr₁₁₅, Leu₁₇₂, Ser₁₇₃, Ile₃₀₃, Leu₅₄₂);</p> <p><i>M2-1 gene</i> (Val₁₈₁);</p> <p><i>M2-2 gene</i> (Asn₂₆);</p> <p><i>N gene</i> (His₂₁₆);</p> <p><i>SH gene</i> (Thr₃₂, Thr₃₇, Thr₄₉);</p> <p><i>L gene</i> (Leu₅₆, Ala₁₇₃, Arg₂₆₀, Val₇₁₅, Ala₁₇₁₂, Asn₁₇₃₆, Ile₂₀₆₉)</p>	17	<p><i>F gene</i> (Arg₄₁₉); <i>NS2 gene</i> (Thr₅); <i>L gene</i> (Ile₁₁₆₆, Arg₂₀₆₆); <i>G gene</i> (Val₂₆₉, Pro₂₇₀, Asp₂₉₄)</p>
9	<p><i>G gene</i> (Thr₁₀₇, Arg₁₃₆, Asn₂₁₃, Leu₂₂₁, Ala₂₂₇, Pro₂₅₅, Asn₃₀₃, Ile₂₇₉, Lys₃₁₂);</p> <p><i>F gene</i> (Ala₁₀₃, Pro₁₂₅, Leu₁₇₂, Ser₁₇₃, Thr₁₉);</p> <p><i>M2-1 gene</i> (Val₁₈₁);</p> <p><i>M2-2 gene</i> (Asn₂₆);</p> <p><i>M gene</i> (Met₇₃);</p> <p><i>N gene</i> (His₂₁₆);</p> <p><i>SH gene</i> (Ile₂₆, Thr₄₉);</p>	18	<p><i>P gene</i> (Ile₆₀);</p> <p><i>L gene</i> (Ile₁₀₅, His₁₄₁, Asn₁₈₃)</p>

	<i>L gene</i> (Leu ₅₆ , Val ₇₁₅ , Ala ₁₇₁₂ , Asn ₁₇₃₆ , Tyr ₁₈₉₇ , Val ₅₁₄)		
		19	<i>P gene</i> (Ile ₆₀); <i>L gene</i> (Ile ₁₀₅ , Asn ₁₈₃); <i>G gene</i> (Pro ₃₁₅)
		20	<i>G gene</i> (His ₉₀ , Phe ₉₁ , Asn ₂₂₅ , Ile ₂₇₃ , Thr ₃₀₁); <i>F gene</i> (Gln ₆₈); <i>NS2 gene</i> (Arg ₈₀); <i>L gene</i> (Ser ₁₈₄)

720

721

722 References

- 723 1. Pneumonia Etiology Research for Child Health Study, G. (2019) 'Causes of severe
724 pneumonia requiring hospital admission in children without HIV infection from Africa
725 and Asia: the PERCH multi-country case-control study', *Lancet*, 394/10200: 757-779.
- 726 2. Scheltema, N. M., Gentile, A., Lucion, F., Nokes, D. J., Munywoki, P. K., Madhi, S. A.,
727 Groome, M. J., Cohen, C., Moyes, J., Thorburn, K., Thamthitawat, S., Oshitani, H.,
728 Lupisan, S. P., Gordon, A., Sanchez, J. F., O'Brien, K. L., Gessner, B. D., Sutanto, A.,
729 Mejias, A., Ramilo, O., Khuri-Bulos, N., Halasa, N., de-Paris, F., Pires, M. R., Spaeder,
730 M. C., Paes, B. A., Simoes, E. A. F., Leung, T. F., da Costa Oliveira, M. T., de Freitas
731 Lazaro Emediato, C. C., Bassat, Q., Butt, W., Chi, H., Aamir, U. B., Ali, A., Lucero, M.
732 G., Fasce, R. A., Lopez, O., Rath, B. A., Polack, F. P., Papenburg, J., Roglic, S., Ito, H.,
733 Goka, E. A., Grobbee, D. E., Nair, H. and Bont, L. J. (2017) 'Global respiratory syncytial
734 virus-associated mortality in young children (RSV GOLD): a retrospective case series',
735 *Lancet Glob Health*, 5/10: e984-e991.
- 736 3. Dowell, S. F., Anderson, L. J., Gary, H. E., Jr., Erdman, D. D., Plouffe, J. F., File, T. M.,
737 Jr., Marston, B. J. and Breiman, R. F. (1996) 'Respiratory syncytial virus is an
738 important cause of community-acquired lower respiratory infection among
739 hospitalized adults', *J Infect Dis*, 174/3: 456-462.
- 740 4. Sullender, W. M. (2000) 'Respiratory syncytial virus genetic and antigenic diversity',
741 *Clin Microbiol Rev*, 13/1: 1-15, table of contents.
- 742 5. Melero, J. A., Garcia-Barreno, B., Martinez, I., Pringle, C. R. and Cane, P. A. (1997)
743 'Antigenic structure, evolution and immunobiology of human respiratory syncytial
744 virus attachment (G) protein', *J Gen Virol*, 78 (Pt 10): 2411-2418.
- 745 6. Trento, A., Casas, I., Calderon, A., Garcia-Garcia, M. L., Calvo, C., Perez-Brena, P. and
746 Melero, J. A. (2010) 'Ten years of global evolution of the human respiratory syncytial
747 virus BA genotype with a 60-nucleotide duplication in the G protein gene', *J Virol*,
748 84/15: 7500-7512.
- 749 7. Zlateva, K. T., Vijgen, L., Dekeersmaecker, N., Naranjo, C. and Van Ranst, M. (2007)
750 'Subgroup Prevalence and Genotype Circulation Patterns of Human Respiratory
751 Syncytial Virus in Belgium during Ten Successive Epidemic Seasons', *J Clin Microbiol*,
752 45/9: 3022.
- 753 8. Agoti, C. N., Otieno, J. R., Gitahi, C. W., Cane, P. A. and Nokes, D. J. (2014) 'Rapid
754 spread and diversification of respiratory syncytial virus genotype ON1, Kenya', *Emerg*
755 *Infect Dis*, 20/6: 950-959.
- 756 9. Otieno, J. R., Kamau, E. M., Agoti, C. N., Lewa, C., Otieno, G., Bett, A., Ngama, M.,
757 Cane, P. A. and Nokes, D. J. (2017) 'Spread and Evolution of Respiratory Syncytial
758 Virus A Genotype ON1, Coastal Kenya, 2010-2015', *Emerg Infect Dis*, 23/2: 264-271.
- 759 10. Bose, M. E., He, J., Shrivastava, S., Nelson, M. I., Bera, J., Halpin, R. A., Town, C. D.,
760 Lorenzi, H. A., Noyola, D. E., Falcone, V., Gerna, G., De Beenhouwer, H., Videla, C.,
761 Kok, T., Venter, M., Williams, J. V. and Henrickson, K. J. (2015) 'Sequencing and
762 analysis of globally obtained human respiratory syncytial virus A and B genomes',
763 *PLoS One*, 10/3: e0120098.
- 764 11. Neuzil, K. M. (2016) 'Progress toward a Respiratory Syncytial Virus Vaccine', *Clinical*
765 *and Vaccine Immunology*, 23/3: 186-188.
- 766 12. Gerretsen, H. E. and Sande, C. J. (2017) 'Development of respiratory syncytial virus
767 (RSV) vaccines for infants', *J Infect*, 74 Suppl 1: S143-s146.

- 768 13. Agoti, C. N., Otieno, J. R., Ngama, M., Mwhuri, A. G., Medley, G. F., Cane, P. A. and
769 Nokes, D. J. (2015) 'Successive Respiratory Syncytial Virus Epidemics in Local
770 Populations Arise from Multiple Variant Introductions, Providing Insights into Virus
771 Persistence', *J Virol*, 89/22: 11630-11642.
- 772 14. Otieno, J. R., Agoti, C. N., Gitahi, C. W., Bett, A., Ngama, M., Medley, G. F., Cane, P. A.
773 and Nokes, D. J. (2016) 'Molecular Evolutionary Dynamics of Respiratory Syncytial
774 Virus Group A in Recurrent Epidemics in Coastal Kenya', *J Virol*, 90/10: 4990-5002.
- 775 15. Otieno, J. R., Kamau, E. M., Oketch, J. W., Ngoi, J. M., Gichuki, A. M., Binter, S.,
776 Otieno, G. P., Ngama, M., Agoti, C. N., Cane, P. A., Kellam, P., Cotten, M., Lemey, P.
777 and Nokes, D. J. (2018) 'Erratum: Whole genome analysis of local Kenyan and global
778 sequences unravels the epidemiological and molecular evolutionary dynamics of RSV
779 genotype ON1 strains', *Virus Evol*, 4/2: vey036.
- 780 16. Nyiro, J. U., Munywoki, P., Kamau, E., Agoti, C., Gichuki, A., Etyang, T., Otieno, G. and
781 Nokes, D. J. (2018) 'Surveillance of respiratory viruses in the outpatient setting in
782 rural coastal Kenya: baseline epidemiological observations', *Wellcome Open Res*, 3:
783 89.
- 784 17. Dudas, G., Carvalho, L. M., Bedford, T., Tatem, A. J., Baele, G., Faria, N. R., Park, D. J.,
785 Ladner, J. T., Arias, A., Asogun, D., Bielejec, F., Caddy, S. L., Cotten, M., D'Ambrozio,
786 J., Dellicour, S., Di Caro, A., Diclaro, J. W., Duraffour, S., Elmore, M. J., Fakoli, L. S.,
787 Faye, O., Gilbert, M. L., Geva, S. M., Gire, S., Gladden-Young, A., Gnrke, A., Goba,
788 A., Grant, D. S., Haagmans, B. L., Hiscox, J. A., Jah, U., Kugelman, J. R., Liu, D., Lu, J.,
789 Malboeuf, C. M., Mate, S., Matthews, D. A., Matranga, C. B., Meredith, L. W., Qu, J.,
790 Quick, J., Pas, S. D., Phan, M. V. T., Pollakis, G., Reusken, C. B., Sanchez-Lockhart, M.,
791 Schaffner, S. F., Schieffelin, J. S., Sealfon, R. S., Simon-Loriere, E., Smits, S. L.,
792 Stoecker, K., Thorne, L., Tobin, E. A., Vandi, M. A., Watson, S. J., West, K., Whitmer,
793 S., Wiley, M. R., Winnicki, S. M., Wohl, S., Wolfel, R., Yozwiak, N. L., Andersen, K. G.,
794 Blyden, S. O., Bolay, F., Carroll, M. W., Dahn, B., Diallo, B., Formenty, P., Fraser, C.,
795 Gao, G. F., Garry, R. F., Goodfellow, I., Gunther, S., Happi, C. T., Holmes, E. C., Kargbo,
796 B., Keita, S., Kellam, P., Koopmans, M. P. G., Kuhn, J. H., Loman, N. J., Magassouba,
797 N., Naidoo, D., Nichol, S. T., Nyenswah, T., Palacios, G., Pybus, O. G., Sabeti, P. C.,
798 Sall, A., Stroher, U., Wurie, I., Suchard, M. A., Lemey, P. and Rambaut, A. (2017)
799 'Virus genomes reveal factors that spread and sustained the Ebola epidemic', *Nature*,
800 544/7650: 309-315.
- 801 18. Faria, N. R., Quick, J., Claro, I. M., Theze, J., de Jesus, J. G., Giovanetti, M., Kraemer,
802 M. U. G., Hill, S. C., Black, A., da Costa, A. C., Franco, L. C., Silva, S. P., Wu, C. H.,
803 Raghwan, J., Cauchemez, S., du Plessis, L., Verotti, M. P., de Oliveira, W. K., Carmo,
804 E. H., Coelho, G. E., Santelli, A., Vinhal, L. C., Henriques, C. M., Simpson, J. T., Loose,
805 M., Andersen, K. G., Grubaugh, N. D., Somasekar, S., Chiu, C. Y., Munoz-Medina, J. E.,
806 Gonzalez-Bonilla, C. R., Arias, C. F., Lewis-Ximenez, L. L., Baylis, S. A., Chieppe, A. O.,
807 Aguiar, S. F., Fernandes, C. A., Lemos, P. S., Nascimento, B. L. S., Monteiro, H. A. O.,
808 Siqueira, I. C., de Queiroz, M. G., de Souza, T. R., Bezerra, J. F., Lemos, M. R., Pereira,
809 G. F., Loudal, D., Moura, L. C., Dhalia, R., Franca, R. F., Magalhaes, T., Marques, E. T.,
810 Jr., Jaenisch, T., Wallau, G. L., de Lima, M. C., Nascimento, V., de Cerqueira, E. M., de
811 Lima, M. M., Mascarenhas, D. L., Neto, J. P. M., Levin, A. S., Tozetto-Mendoza, T. R.,
812 Fonseca, S. N., Mendes-Correa, M. C., Milagres, F. P., Segurado, A., Holmes, E. C.,
813 Rambaut, A., Bedford, T., Nunes, M. R. T., Sabino, E. C., Alcantara, L. C. J., Loman, N.

- 814 J. and Pybus, O. G. (2017) 'Establishment and cryptic transmission of Zika virus in
815 Brazil and the Americas', *Nature*, 546/7658: 406-410.
- 816 19. Lemey, P., Rambaut, A., Bedford, T., Faria, N., Bielejec, F., Baele, G., Russell, C. A.,
817 Smith, D. J., Pybus, O. G., Brockmann, D. and Suchard, M. A. (2014) 'Unifying viral
818 genetics and human transportation data to predict the global transmission dynamics
819 of human influenza H3N2', *PLoS Pathog*, 10/2: e1003932.
- 820 20. Faria, N. R., Suchard, M. A., Rambaut, A., Streicker, D. G. and Lemey, P. (2013)
821 'Simultaneously reconstructing viral cross-species transmission history and
822 identifying the underlying constraints', *Philos Trans R Soc Lond B Biol Sci*, 368/1614:
823 20120196.
- 824 21. Zehender, G., Veo, C., Ebranati, E., Carta, V., Rovida, F., Percivalle, E., Moreno, A.,
825 Lelli, D., Calzolari, M., Lavazza, A., Chiapponi, C., Baioni, L., Capelli, G., Ravagnan, S.,
826 Da Rold, G., Lavezzo, E., Palu, G., Baldanti, F., Barzon, L. and Galli, M. (2017)
827 'Reconstructing the recent West Nile virus lineage 2 epidemic in Europe and Italy
828 using discrete and continuous phylogeography', *PLoS One*, 12/7: e0179679.
- 829 22. Nokes, D. J., Ngama, M., Bett, A., Abwao, J., Munywoki, P., English, M., Scott, J. A.,
830 Cane, P. A. and Medley, G. F. (2009) 'Incidence and severity of respiratory syncytial
831 virus pneumonia in rural Kenyan children identified through hospital surveillance',
832 *Clin Infect Dis*, 49/9: 1341-1349.
- 833 23. Scott, J. A. G., Bauni, E., Moisi, J. C., Ojal, J., Gatakaa, H., Nyundo, C., Molyneux, C. S.,
834 Kombe, F., Tsofa, B., Marsh, K., Peshu, N. and Williams, T. N. (2012) 'Profile: The Kilifi
835 Health and Demographic Surveillance System (KHDSS)', *International Journal of
836 Epidemiology*, 41/3: 650-657.
- 837 24. Hammitt, L. L., Kazungu, S., Welch, S., Bett, A., Onyango, C. O., Gunson, R. N., Scott,
838 J. A. and Nokes, D. J. (2011) 'Added value of an oropharyngeal swab in detection of
839 viruses in children hospitalized with lower respiratory tract infection', *J Clin
840 Microbiol*, 49/6: 2318-2320.
- 841 25. Kamau, E., Agoti, C. N., Lewa, C. S., Oketch, J., Owor, B. E., Otieno, G. P., Bett, A.,
842 Cane, P. A. and Nokes, D. J. (2017) 'Recent sequence variation in probe binding site
843 affected detection of respiratory syncytial virus group B by real-time RT-PCR', *J Clin
844 Virol*, 88: 21-25.
- 845 26. Agoti, C. N., Otieno, J. R., Munywoki, P. K., Mwihuri, A. G., Cane, P. A., Nokes, D. J.,
846 Kellam, P. and Cotten, M. (2015) 'Local evolutionary patterns of human respiratory
847 syncytial virus derived from whole-genome sequencing', *J Virol*, 89/7: 3444-3454.
- 848 27. Langmead, B. and Salzberg, S. L. (2012) 'Fast gapped-read alignment with Bowtie 2',
849 *Nature Methods*, 9: 357.
- 850 28. Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G.,
851 Abecasis, G. and Durbin, R. (2009) 'The Sequence Alignment/Map format and
852 SAMtools', *Bioinformatics*, 25/16: 2078-2079.
- 853 29. Quinlan, A. R. and Hall, I. M. (2010) 'BEDTools: a flexible suite of utilities for
854 comparing genomic features', *Bioinformatics*, 26/6: 841-842.
- 855 30. Katoh, K. and Standley, D. M. (2013) 'MAFFT Multiple Sequence Alignment Software
856 Version 7: Improvements in Performance and Usability', *Molecular Biology and
857 Evolution*, 30/4: 772-780.
- 858 31. Stamatakis, A. (2014) 'RAxML version 8: a tool for phylogenetic analysis and post-
859 analysis of large phylogenies', *Bioinformatics*, 30/9: 1312-1313.

- 860 32. Rambaut, A., Lam, T. T., Max Carvalho, L. and Pybus, O. G. (2016) 'Exploring the
861 temporal structure of heterochronous sequences using TempEst (formerly Path-O-
862 Gen)', *Virus Evolution*, 2/1: vew007-vew007.
- 863 33. Posada, D. (2008) 'jModelTest: phylogenetic model averaging', *Mol Biol Evol*, 25/7:
864 1253-1256.
- 865 34. Suchard, M. A., Lemey, P., Baele, G., Ayres, D. L., Drummond, A. J. and Rambaut, A.
866 (2018) 'Bayesian phylogenetic and phylodynamic data integration using BEAST 1.10',
867 *Virus Evolution*, 4/1.
- 868 35. Baele, G., Li, W. L., Drummond, A. J., Suchard, M. A. and Lemey, P. (2013) 'Accurate
869 model selection of relaxed molecular clocks in bayesian phylogenetics', *Mol Biol Evol*,
870 30/2: 239-243.
- 871 36. Lemey, P., Rambaut, A., Drummond, A. J. and Suchard, M. A. (2009) 'Bayesian
872 phylogeography finds its roots', *PLoS Comput Biol*, 5/9: e1000520.
- 873 37. Faria, N. R., Hodges-Mameletzis, I., Silva, J. C., Rodes, B., Erasmus, S., Paolucci, S.,
874 Ruelle, J., Pieniazek, D., Taveira, N., Trevino, A., Goncalves, M. F., Jallow, S., Xu, L.,
875 Camacho, R. J., Soriano, V., Goubau, P., de Sousa, J. D., Vandamme, A. M., Suchard,
876 M. A. and Lemey, P. (2012) 'Phylogeographical footprint of colonial history in the
877 global dispersal of human immunodeficiency virus type 2 group A', *J Gen Virol*, 93/Pt
878 4: 889-899.
- 879 38. Bielejec, F., Baele, G., Vrancken, B., Suchard, M. A., Rambaut, A. and Lemey, P.
880 (2016) 'Spread3: Interactive Visualization of Spatiotemporal History and Trait
881 Evolutionary Processes', *Mol Biol Evol*, 33/8: 2167-2169.
- 882 39. Weaver, S., Shank, S. D., Spielman, S. J., Li, M., Muse, S. V. and Kosakovsky Pond, S. L.
883 (2018) 'Datamonkey 2.0: A Modern Web Application for Characterizing Selective and
884 Other Evolutionary Processes', *Mol Biol Evol*, 35/3: 773-777.
- 885 40. Nguyen, L. T., Schmidt, H. A., von Haeseler, A. and Minh, B. Q. (2015) 'IQ-TREE: a fast
886 and effective stochastic algorithm for estimating maximum-likelihood phylogenies',
887 *Mol Biol Evol*, 32/1: 268-274.
- 888 41. Parker, J., Rambaut, A. and Pybus, O. G. (2008) 'Correlating viral phenotypes with
889 phylogeny: accounting for phylogenetic uncertainty', *Infect Genet Evol*, 8/3: 239-246.
- 890 42. Sagulenko, P., Puller, V. and Neher, R. A. (2018) 'TreeTime: Maximum-likelihood
891 phylodynamic analysis', *Virus Evol*, 4/1: vex042.
- 892 43. Botosso, V. F., Zanotto, P. M., Ueda, M., Arruda, E., Gilio, A. E., Vieira, S. E., Stewien,
893 K. E., Peret, T. C., Jamal, L. F., Pardini, M. I., Pinho, J. R., Massad, E., Sant'anna, O. A.,
894 Holmes, E. C., Durigon, E. L. and Consortium, V. (2009) 'Positive selection results in
895 frequent reversible amino acid replacements in the G protein gene of human
896 respiratory syncytial virus', *PLoS Pathog*, 5/1: e1000254.
- 897 44. Bin, L., Liu, H., Tabor, D. E., Tovchigrechko, A., Qi, Y., Ruzin, A., Esser, M. T. and Jin, H.
898 (2019) 'Emergence of new antigenic epitopes in the glycoproteins of human
899 respiratory syncytial virus collected from a US surveillance study, 2015-17', *Sci Rep*,
900 9/1: 3898.
- 901 45. Duffy, S., Shackelton, L. A. and Holmes, E. C. (2008) 'Rates of evolutionary change in
902 viruses: patterns and determinants', *Nat Rev Genet*, 9/4: 267-276.
- 903 46. Aldous, J. L., Pond, S. K., Poon, A., Jain, S., Qin, H., Kahn, J. S., Kitahata, M., Rodriguez,
904 B., Dennis, A. M., Boswell, S. L., Haubrich, R. and Smith, D. M. (2012) 'Characterizing
905 HIV transmission networks across the United States', *Clin Infect Dis*, 55/8: 1135-
906 1143.

- 907 47. Smith, D. J., Lapedes, A. S., de Jong, J. C., Bestebroer, T. M., Rimmelzwaan, G. F.,
908 Osterhaus, A. D. and Fouchier, R. A. (2004) 'Mapping the antigenic and genetic
909 evolution of influenza virus', *Science*, 305/5682: 371-376.
- 910 48. Zhu, Q., Lu, B., McTamney, P., Palaszynski, S., Diallo, S., Ren, K., Ulbrandt, N. D.,
911 Kallewaard, N., Wang, W., Fernandes, F., Wong, S., Svabek, C., Moldt, B., Esser, M. T.,
912 Jing, H. and Suzich, J. A. (2018) 'Prevalence and Significance of Substitutions in the
913 Fusion Protein of Respiratory Syncytial Virus Resulting in Neutralization Escape From
914 Antibody MEDI8897', *J Infect Dis*, 218/4: 572-580.
- 915 49. Domachowske, J. B., Khan, A. A., Esser, M. T., Jensen, K., Takas, T., Villafana, T.,
916 Dubovsky, F. and Griffin, M. P. (2018) 'Safety, Tolerability and Pharmacokinetics of
917 MEDI8897, an Extended Half-life Single-dose Respiratory Syncytial Virus Prefusion F-
918 targeting Monoclonal Antibody Administered as a Single Dose to Healthy Preterm
919 Infants', *Pediatr Infect Dis J*, 37/9: 886-892.
- 920 50. Zhu, Q., McLellan, J. S., Kallewaard, N. L., Ulbrandt, N. D., Palaszynski, S., Zhang, J.,
921 Moldt, B., Khan, A., Svabek, C., McAuliffe, J. M., Wrapp, D., Patel, N. K., Cook, K. E.,
922 Richter, B. W. M., Ryan, P. C., Yuan, A. Q. and Suzich, J. A. (2017) 'A highly potent
923 extended half-life antibody as a potential RSV vaccine surrogate for all infants', *Sci*
924 *Transl Med*, 9/388.
- 925 51. Kosakovsky Pond, S. L. and Frost, S. D. (2005) 'Not so different after all: a comparison
926 of methods for detecting amino acid sites under selection', *Mol Biol Evol*, 22/5: 1208-
927 1222.
- 928 52. Chen, X., Xu, B., Guo, J., Li, C., An, S., Zhou, Y., Chen, A., Deng, L., Fu, Z., Zhu, Y., Liu,
929 C., Xu, L., Wang, W., Shen, K. and Xie, Z. (2018) 'Genetic variations in the fusion
930 protein of respiratory syncytial virus isolated from children hospitalized with
931 community-acquired pneumonia in China', *Sci Rep*, 8/1: 4491.
932

Figure 1 (A)

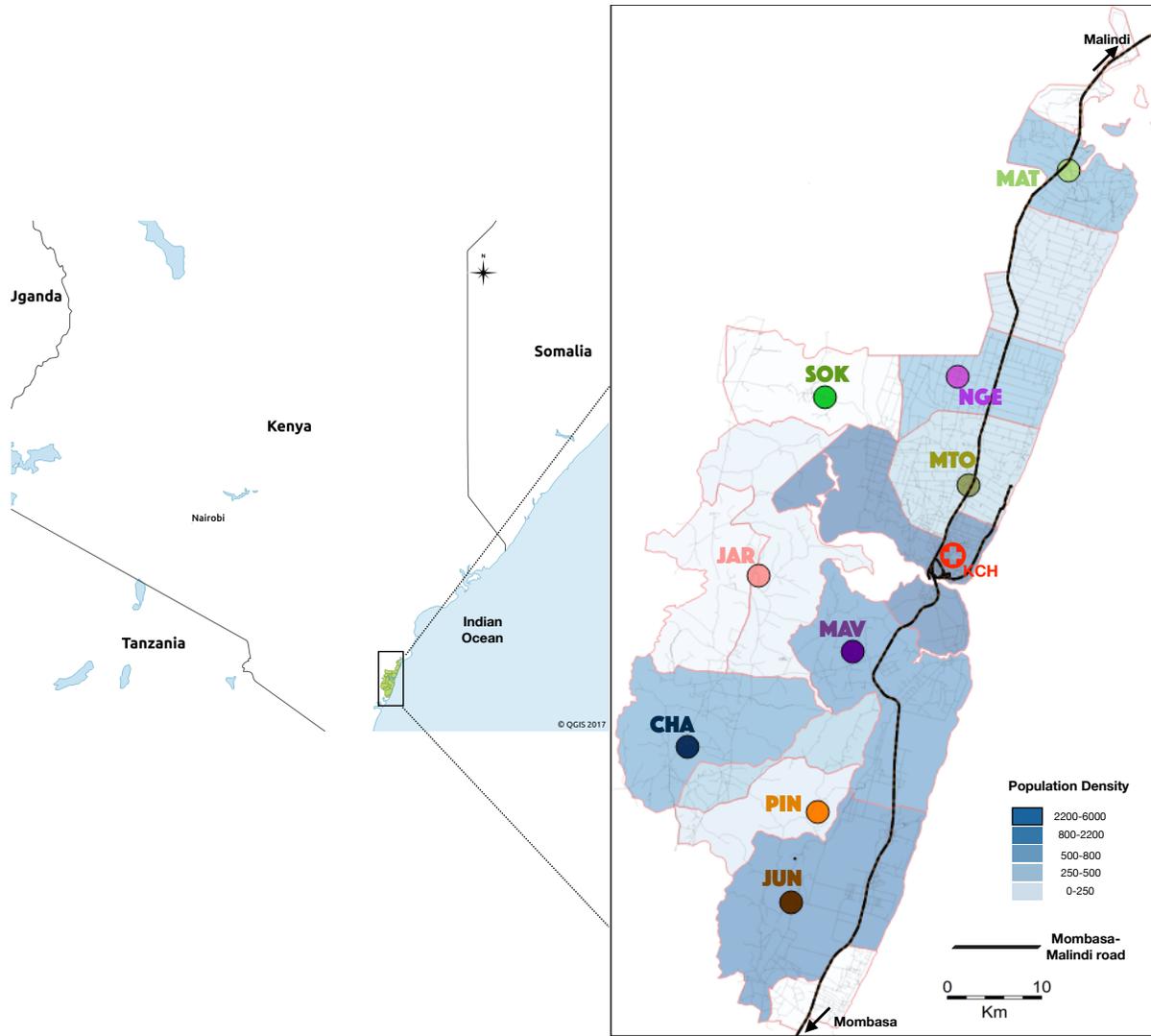


Figure 2 (A)

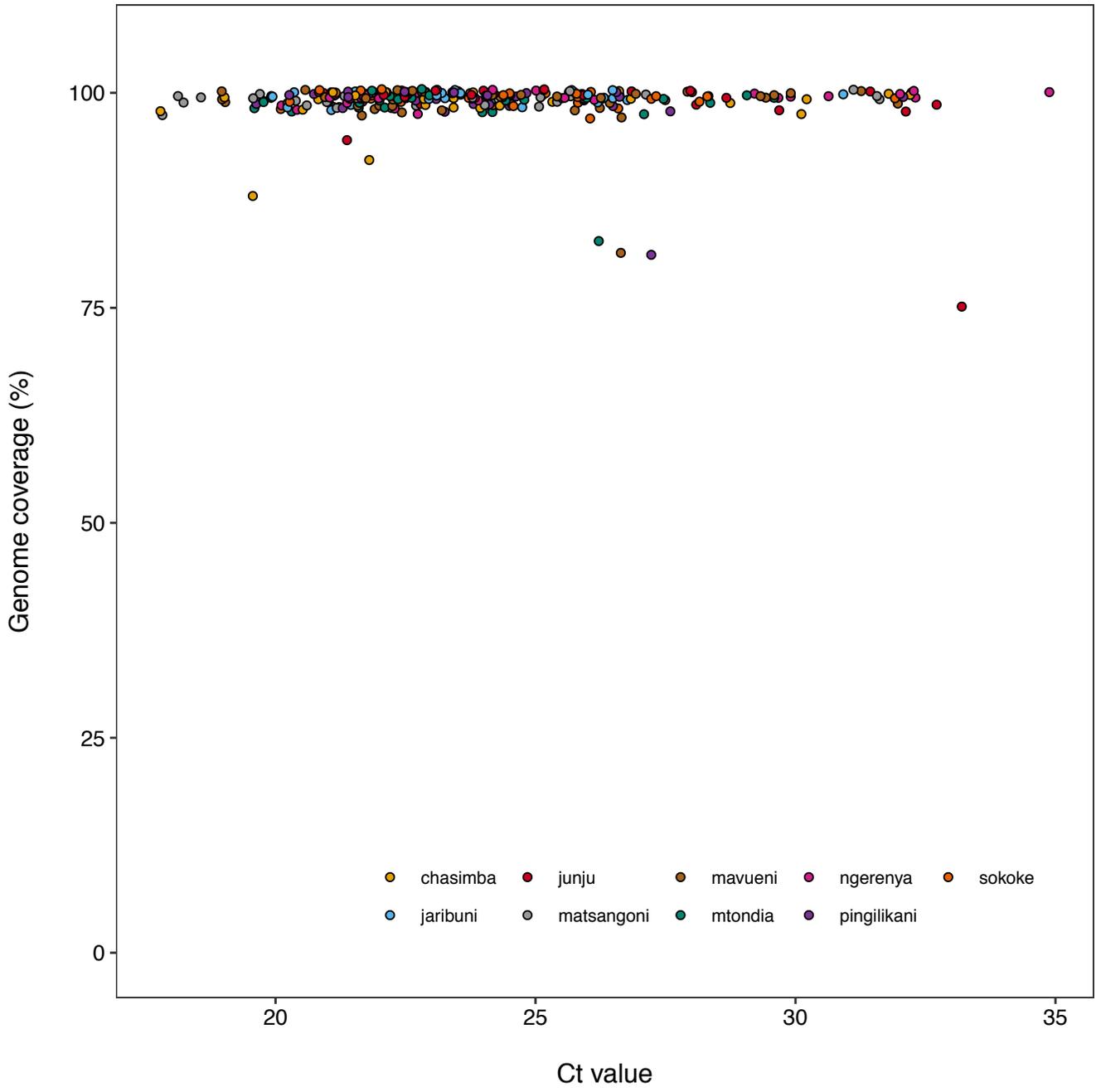


Figure 2(C)

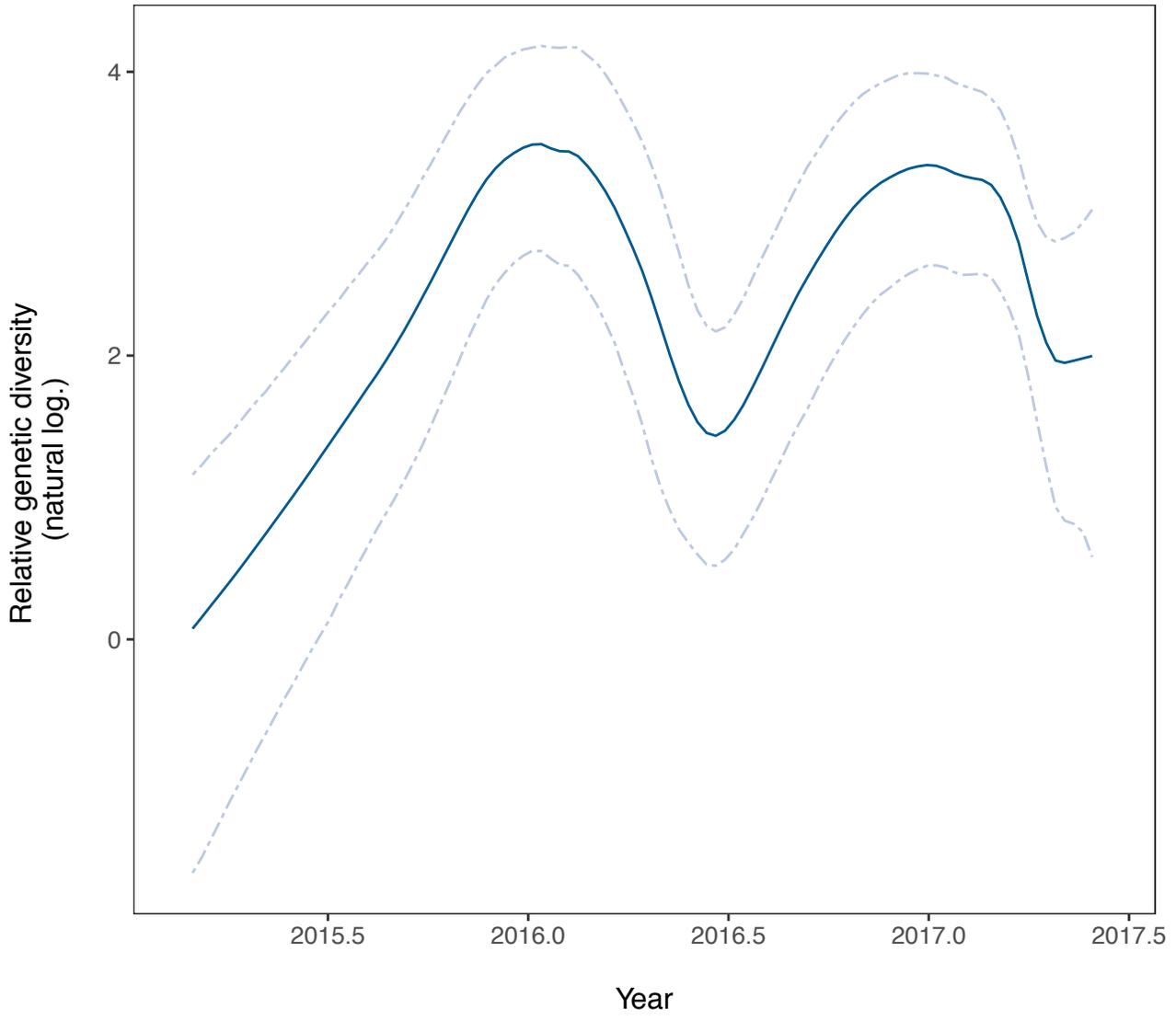


Figure 2 (D)

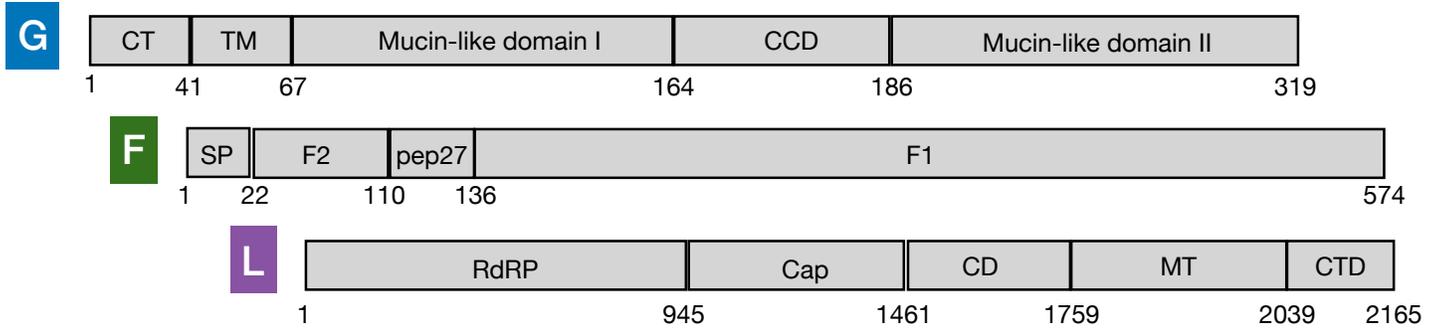
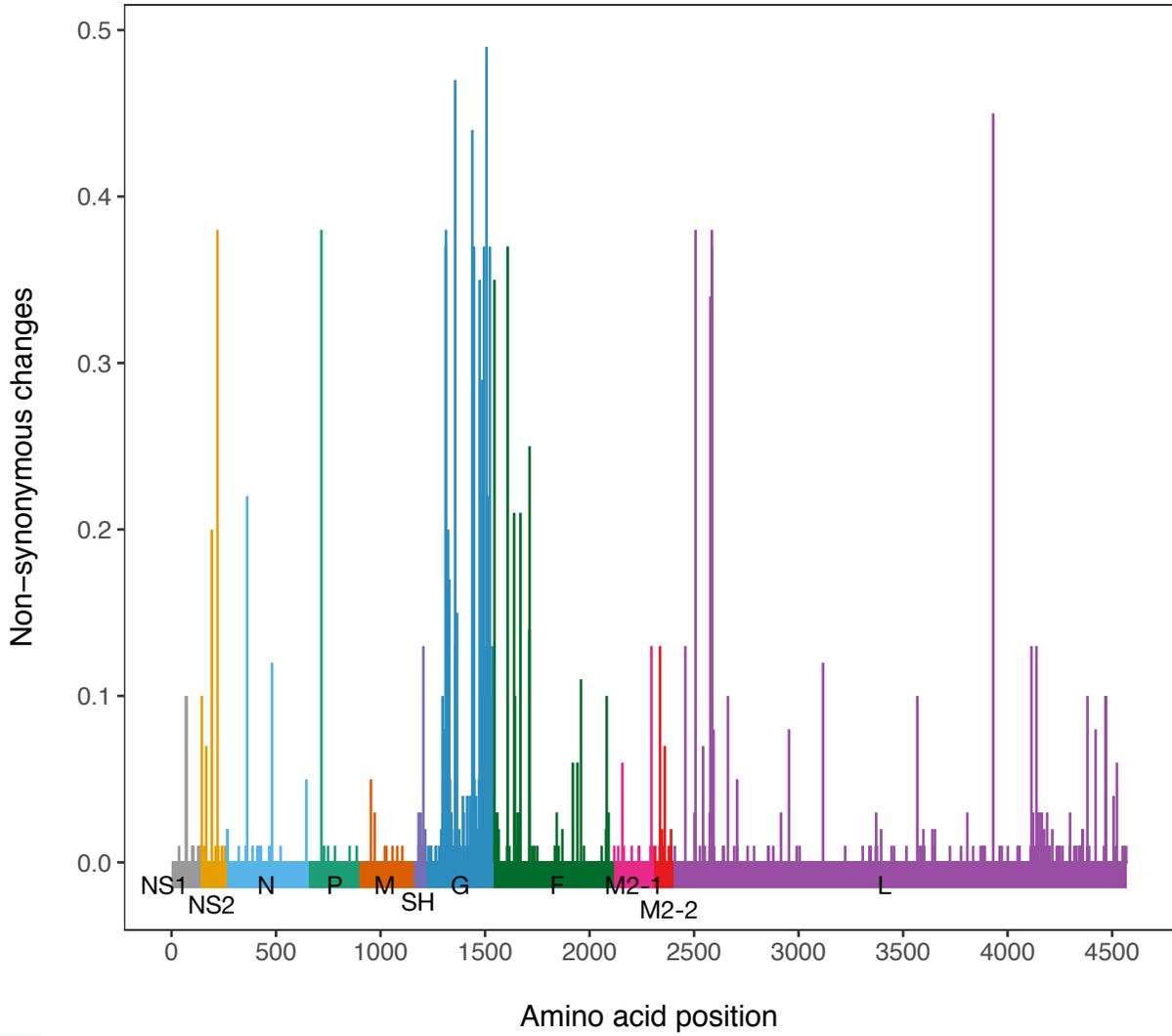


Figure 3(A)



	I	II	III	IV	V	VI	VII	VIII
CHA			⊗		⊗	⊗	⊗	⊗
JAR		⊗	⊗			⊗	⊗	⊗
JUN		⊗	⊗		⊗	⊗	⊗	⊗
MAV		⊗	⊗	⊗			⊗	⊗
MAT			⊗		⊗		⊗	⊗
PIN		⊗	⊗	⊗	⊗	⊗	⊗	⊗
NGE			⊗		⊗		⊗	⊗
SOK	⊗	⊗	⊗	⊗	⊗			
MTO	⊗	⊗	⊗		⊗	⊗	⊗	⊗

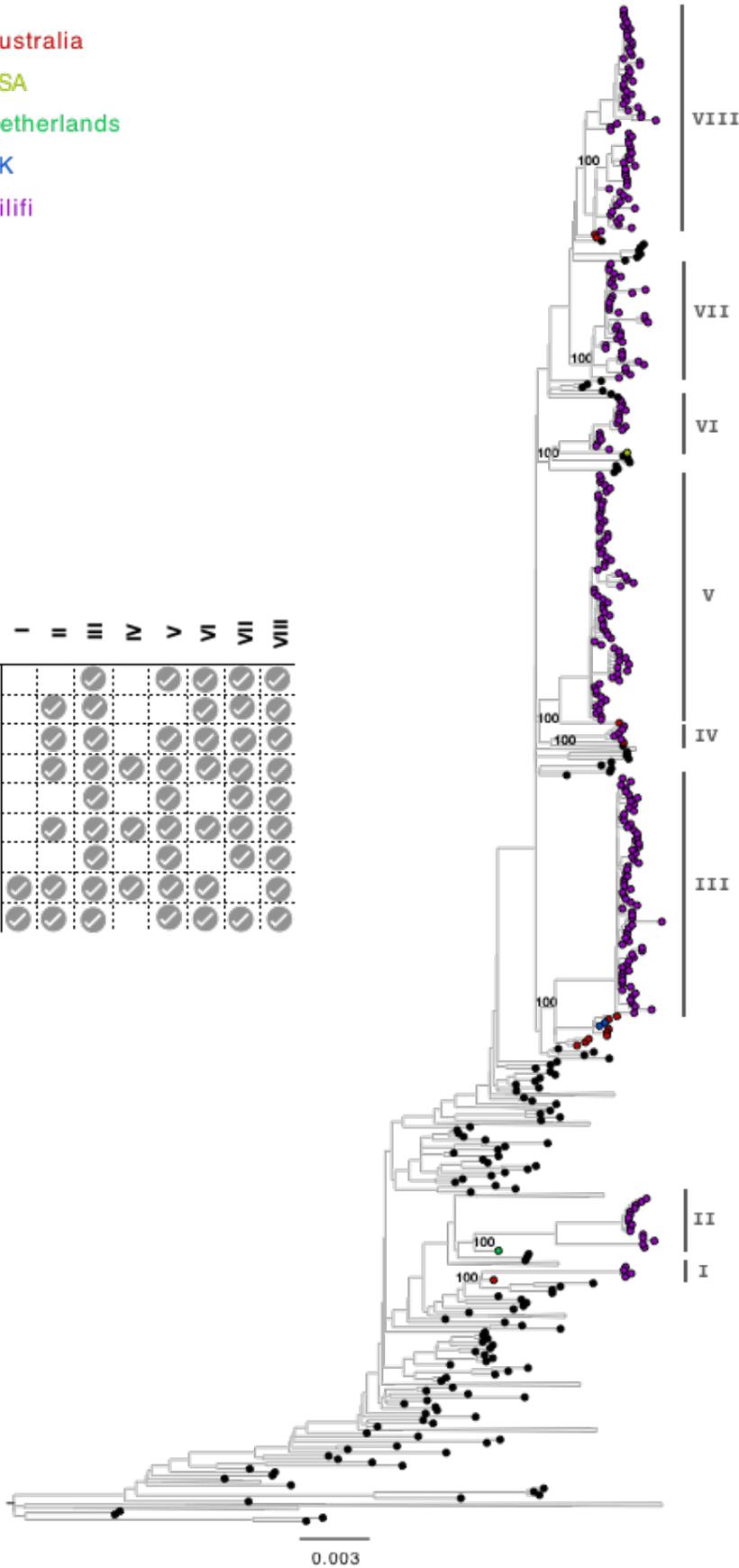


Figure 3(B)

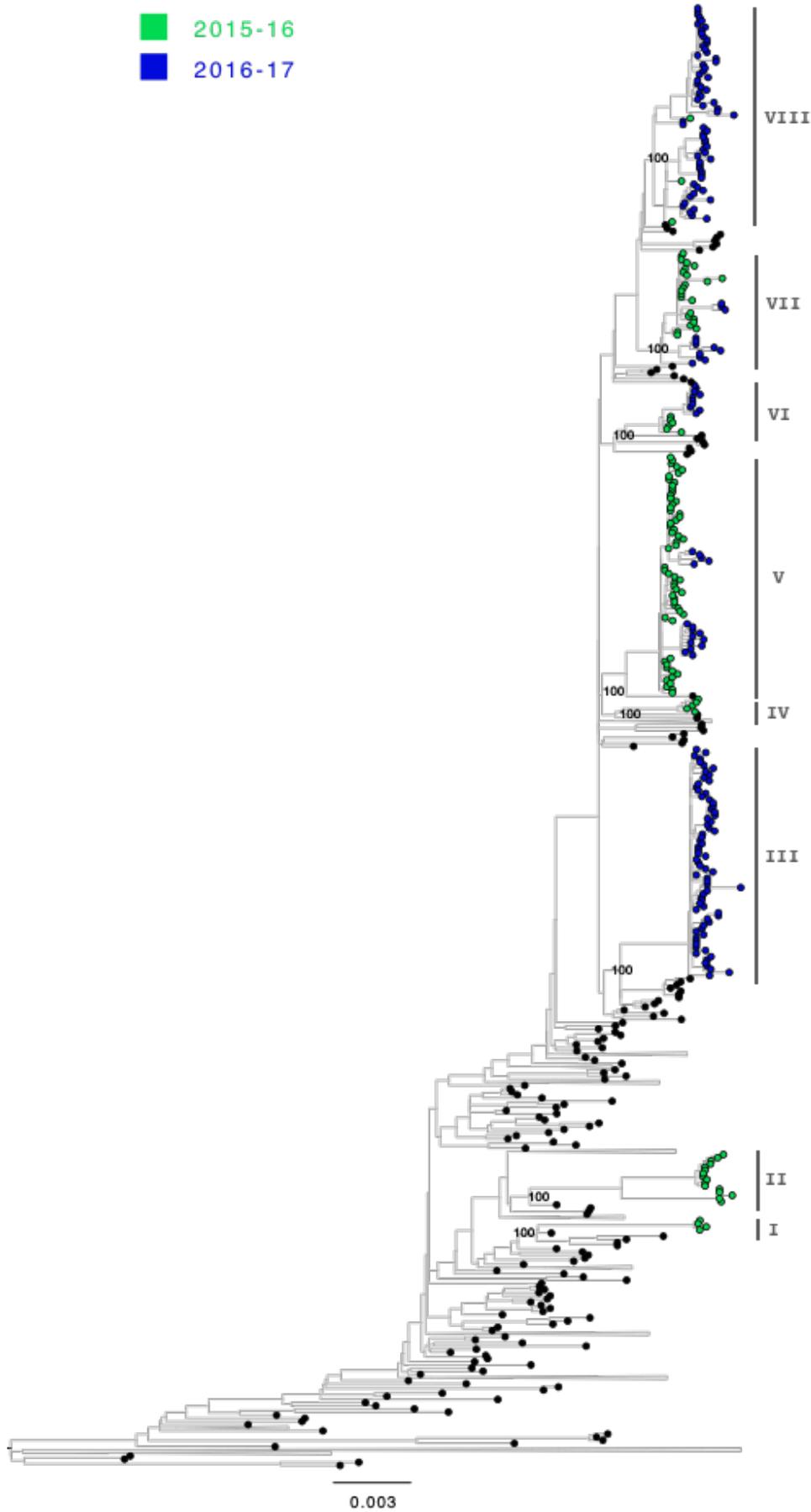


Figure 4(B)

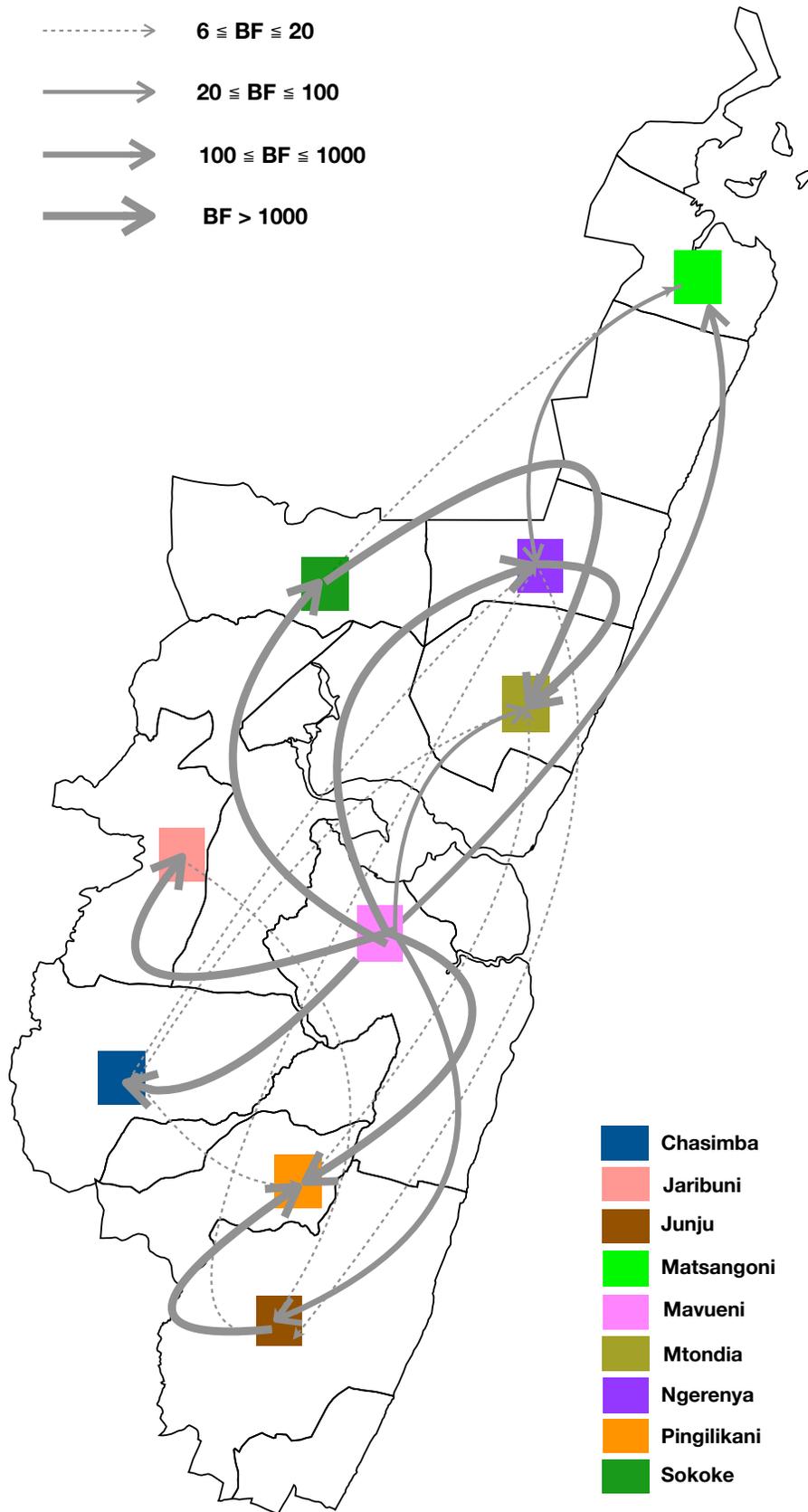


Figure 5(A)

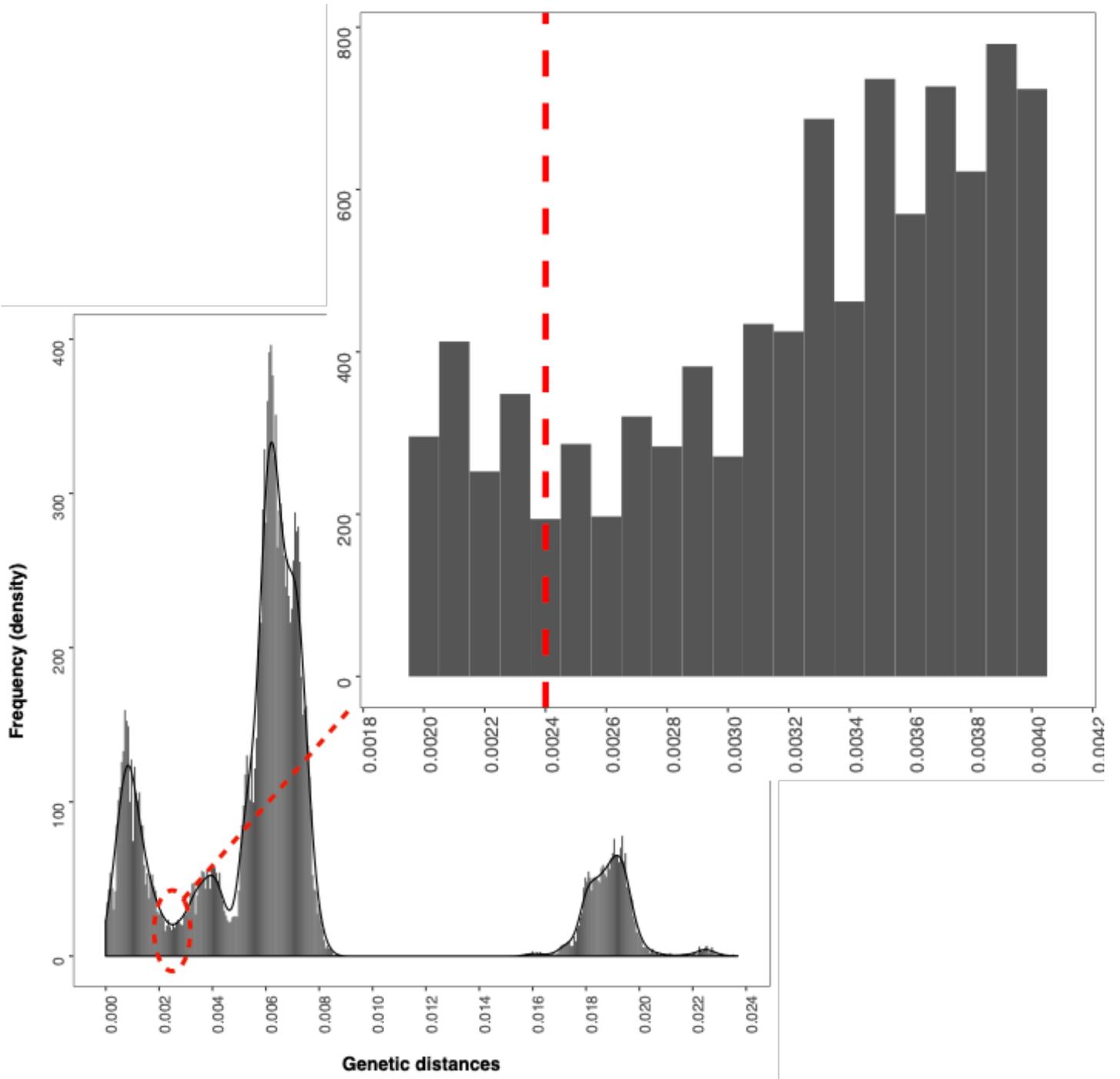


Figure 5(B)

2015/16		2016/17	
1	0.0013362	10	0.0018424
2	0.0021504	11	0.0003305
3	0.0002829	12	1E-06
4	0.0022764	13	0.0019125
5	0.0012214	14	0.0018402
6	0.0012025	15	0.0021738
7	0.0009352	16	0.000208
8	0.0006647	17	0.0010525
9	0.0005339	18	0.0015105
		19	0.0017768
		20	0.0019813

Figure 6(A)

Cluster	1	2	3	4	5	6	7	8	9
Cha									
Jar									
Jun									
Mat									
Mav									
Mto									
Nge									
Pin									
Sok									
# viruses	60	5	21	3	3	6	11	5	4

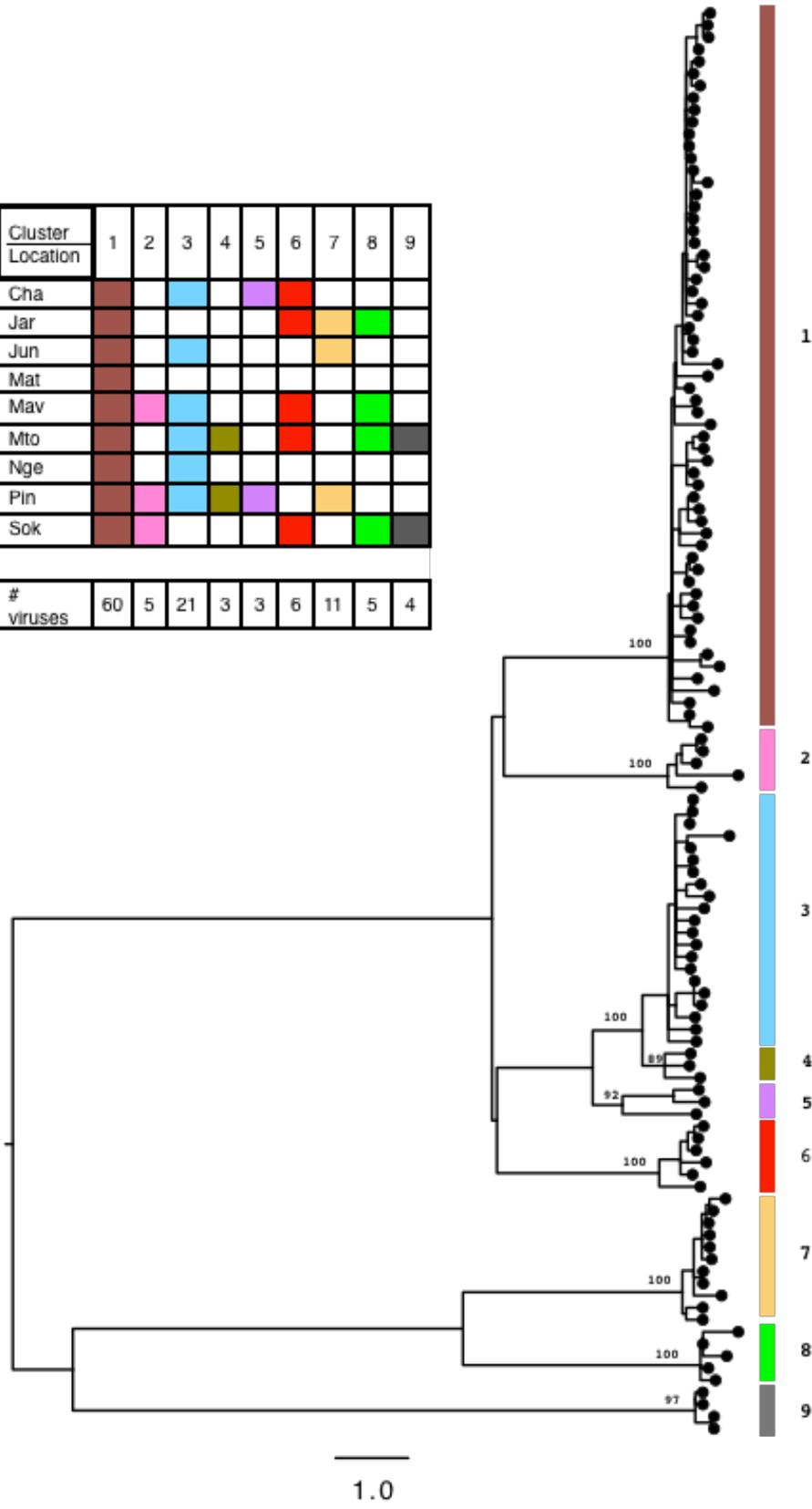


Figure 6(B)

