1 **Title:**

2 Surveillance of bat coronaviruses in Kenya identifies relatives of human coronaviruses NL63

3 and 229E and their recombination history

4

5 **Running title:**

6 Bat origin of human coronaviruses

7

8 **Authors:**

9 Ying Tao[1#], Mang Shi[2#], Christina Chommanard[1], Krista Queen[1], Jing Zhang[1], Wanda

10 Markotter[3], Ivan V. Kuzmin[4]†, Edward C. Holmes[2], Suxiang Tong[1*]

11

12 **Affiliations:**

13 [1] Division of Viral Diseases, Centers for Disease Control and Prevention, Atlanta, GA 30333,

14 USA; [2]Marie Bashir Institute for Infectious Diseases and Biosecurity, Charles Perkins Centre,

15 School of Life and Environmental Sciences and Sydney Medical School, The University of

16 Sydney, Sydney, Australia; [3]Centre for Viral Zoonoses, Department of Medical Virology,

17 Faculty of Health Sciences, University of Pretoria, Pretoria, South Africa; [4]Division of High

18 Consequence Pathogens and Pathology, Centers for Disease Control and Prevention, Atlanta,

19 GA 30333, USA.

20 # Y.T. and M.S. contributed equally to this work

21 † Present address: Department of Pathology, University of Texas Medical Branch, Galveston,
22 TX 77555, USA.

23 * Correspondence to:  Dr. Suxiang Tong,        11600 Clifton Rd, mail stop G18, CDC,

24 Atlanta, GA 30333; Tel: 4046391372; Email: sot1@cdc.gov.

25 The findings and conclusions in this report are those of the author(s) and do not necessarily
26 represent the official position of the Centers for Disease Control and Prevention.
27

28 **Type of Publication:** 'Full length' paper

29 **Word count:** Abstract (155); Importance (106); Text body (4618)

30   **ABSTRACT**

31   Bats harbor a large diversity of coronaviruses (CoVs), several of which are related to

32   zoonotic pathogens that cause severe disease in humans. Our screening of bat samples

33   collected in Kenya during 2007-2010 not only detected RNA from several novel CoVs but,

34   more significantly, identified sequences that were closely related to human CoVs NL63 and

35   229E, suggesting that these two human viruses originate from bats. We also demonstrated

36   that human CoV NL63 is a recombinant between NL63-like viruses circulating in *Triaenops*

37   bats and 229E-like viruses circulating in *Hipposideros* bats, with the break-point located near

38   5' and 3' end of the spike (S) protein gene. In addition, two further inter-species

39   recombination events involving the S gene were identified, suggesting that this region may

40   represent a recombination "hotspot" in CoV genomes. Finally, using a combination of

41   phylogenetic and distance-based approaches we showed that genetic diversity of bat CoVs is

42   primarily structured by host species and subsequently by geographic distances.

43

44   **IMPORTANCE**

45   Understanding the driving forces of cross-species virus transmission is central to

46   understanding the nature of disease emergence. Previous studies have demonstrated that bats

47   are the ultimate reservoir hosts for a number of coronaviruses (CoVs) including ancestors of

48   SARS-CoV, MERS-CoV, and HCoV-229E. However, the evolutionary pathways of bat

49   CoVs remain elusive. We provide evidence for natural recombination between distantly-

50   related African bat coronaviruses associated with *Triaenops afer* and *Hipposideros* sp. bats

51   that resulted in a NL-63 like virus, an ancestor of the human pathogen HCoV-NL63. These

52   results suggest that inter-species recombination may play an important role in CoV evolution

53   and the emergence of novel CoVs with zoonotic potential.

54 **INTRODUCTION**

55 Coronaviruses (CoVs) (subfamily *Coronavirinae*, family *Coronaviridae*, order Nidovirales)

56 are common infectious agents that infect a wide range of hosts including humans, causing

57 respiratory, gastrointestinal, liver, and neurologic diseases, and that possess the largest

58 genomes of any RNA viruses described to date (1). The subfamily *Coronavirinae* is currently

59 classified into four genera: *Alphacoronavirus, Betacoronavirus*, *Gammacoronavirus*, and

60 *Deltacoronavirus* (2). The alphacoronaviruses (alpha-CoV) and betacoronaviruses (beta-CoV)

61 are exclusively found in mammals while the gammacoronaviruses (gamma-CoV) and

62 deltacoronaviruses (delta-CoV) are mainly associated with birds. Presently, the greatest

63 diversity of alpha- and beta-CoVs has been documented in bats, which in part reflects the

64 more intensive surveillance of these animals since *Rhinolophus* spp. bats were implicated as

65 the reservoir hosts for SARS-related CoVs (3, 4). This surveillance resulted in the discovery

66 of a potential reservoir host (bat) species for another two human CoVs: Human CoV 229E

67 (HCoV-229E), a relative of which is present in *Hipposideros* bats (5, 6), and Middle East

68 respiratory syndrome coronavirus (MERS-CoV), for which related viruses are present in

69 *Pipistrellus*, *Tylonycteris*, and *Neoromicia* bats (7-10), although the most likely reservoir host

70 of human MERS-CoV identified to date is the dromedary camel (11). Most recently HCoV-

71 229E-like CoVs were also identified in camels, although their role in human infection is

72 unknown (12).

73     Africa is a major hotspot of zoonotic emerging diseases. With its rich biodiversity,

74 Africa is inhabited by many bats of different species including those that serve as reservoirs

75 of important zoonotic diseases such as Marburg hemorrhagic fever and rabies (13). Our initial

76 screening demonstrated the presence of diverse CoVs in African bats, including those

77 collected in the southern parts of Kenya during 2006 (14, 15), and in other countries

78 including South Africa, Nigeria, and Ghana (16). Furthermore, recent studies have provided

3

79    strong evidence that HCoV-229E originated from bat viruses circulating in Africa (5),

80    underscoring the zoonotic potential of bat-borne CoVs from this continent.

81        One human coronavirus, HCoV-NL63, was first isolated in 2004 from the aspirate of

82    a 8-month-old boy suffering from pneumonia in the Netherlands (17). While the clinical

83    significance of this virus is debated, it has a worldwide distribution and is known to infect

84    both the upper and lower respiratory tract (18). Based on a phylogeny of the RNA-dependent

85    RNA polymerase (RdRp), HCoV-NL63 is related to another human virus HCoV-229E and

86    had no close relatives identified in bats (16). Although Huynh et al. (19) suggested that a

87    virus (ARCoV.2/2010/USA) isolated from the American tricolored bat (*Perimyotis subflavus*)

88    may share common ancestry with HCoV-NL63, the genetic distance between the two viruses

89    is large, and their close relationship has not been corroborated in other phylogenetic analyses

90    (16, 20).  Nevertheless, the successful passage of HCoV-NL63 in an immortalized bat cell

91    line suggests its potential association with bats (19).

92        As is well appreciated, recombination leads to rapid changes of genetic diversity in

93    RNA viruses (21). CoVs represent a classic example of viruses with high frequencies of

94    homologous recombination through discontinuous RNA synthesis (22). Indeed, under

95    experimental conditions, the recombination frequency can be as high as 25% for the entire

96    CoV genome (23). Recombination in CoVs is also frequently reported under natural

97    conditions, including some emerging human pathogens such as SARS-CoV (24, 25), MERS-

98    CoV (11), HCoV-OC43 (26), and HCoV-NL63 (27), although most reports are between

99    closely related viruses.

100    The Global Disease Detection Program (GDD) of the Centers for Disease Control and

101    Prevention (CDC, Atlanta, GA) is focused on the detection of emerging infectious agents

102    worldwide. One of the GDD projects was directed toward the detection of such potential

103    zoonotic pathogens in African bats. Since the initial study performed during 2006 in Kenya

4

104 (14, 15), an expanded surveillance of bat CoVs has been performed in the same and other

105 countries including Kenya, Nigeria, Democratic Republic of Georgia, Democratic Republic

106 of Congo, Guatemala, and Peru. The project included more bat species and geographic

107 locations, allowing a more thorough investigation of the genetic diversity and ecological

108 dynamics of CoVs circulation in bats. In this study, we performed an ecological and

109 evolutionary characterization of CoVs circulating in Kenya and identified distinct CoVs from

110 *Triaenops afer* and *Hipposideros* sp. bats that are phylogenetically related to HCoV-NL63 in

111 different parts of the genome. Based on this data, we propose a scenario for the origin and

112 evolutionary history of HCoV-NL63 and related viruses.

113

114 **MATERIALS AND METHODS**

115 **Sample collection.** Between 2007 and 2010 a total of 2050 bat specimens were collected

116 from 30 different locations in Kenya (Table S1) in collaboration with the CDC GDD regional

117 country office in Kenya and National Museums of Kenya. The bats were captured using mist-

118 nets, hand nets or manually. The protocol (2096FRAMULX-A3) was approved by the CDC

119 IACUC and by Kenya Wildlife Services. Upon capture, each bat was measured, sexed and

120 identified to species by a trained field biologist. Subsequently, fecal and oral swabs (if

121 possible) were collected in compliance with field protocol and were then transported on dry

122 ice from the field to -80°C storage before further processing.

123

124 **CoV RNA detection.** Each fecal and oral swab was suspended in 200 μL of a phosphate

125 buffered saline. Viral total nucleic acids (TNA) were extracted using the QIAamp Mini Viral

126 Spin kit (Qiagen, Valencia, CA, USA) according to the manufacturer's instructions, followed

127 by semi-nested RT-PCR (SuperScript III One-Step RT-PCR kit and Platinum Taq kit,

128 Invitrogen, San Diego, CA, USA) using primer sets designed to target the conserved genome

129 region of alpha-, beta-, gamma- and delta-CoVs, respectively (15). PCR products of the

130 expected size (~ 400 nucleotides) were purified by gel extraction using the QIAquick Gel

131 Extraction kit (Qiagen, Valencia, CA, USA) and sequenced in both directions on an ABI

132 Prism 3130 automated sequencer (Applied Biosystems, Foster City, CA, USA). As

133 validation, the RT-PCR procedure was repeated for each of the CoV positive specimens.

134

135 **Bat mitochondrial gene sequencing.** Bat species were further confirmed by sequencing the

136 host mitochondrial cytochrome b (cytB) gene in each of the CoV-positive specimens. Both

137 the method and the primers used have been described previously, and a final 1104 bp

138 fragment of the cytB gene was amplified and sequenced as described previously (14, 15).

6

139

140 **Phylogenetic analyses**. This study generated a total of 240 CoV RdRP sequences (402 bp)

141 from Kenyan bats. These sequences were first aligned in MAFFT v7.013 (28), using amino

142 acid sequences as a guide for the nucleotide sequence alignment. Phylogenetic trees were

143 then inferred using the maximum likelihood (ML) method available in PhyML version 3.0

144 (29) assuming a general time-reversible (GTR) model with a discrete gamma distributed rate

145 variation among sites ($\Gamma_4$) and the SPR branch-swapping algorithm. To produce a more

146 condensed data set, we clustered the highly similar sequences from the same geographic

147 location and host species, and randomly chose one or two to represent each cluster. This

148 condensed data set was subsequently combined with 121 reference sequences representative

149 of the genetic diversity of alpha- and beta-CoVs on a global scale taken from GenBank. ML

150 phylogenetic trees of these final alignments were inferred using the same procedure and

151 substitution models as described above.

152

153 **Comparisons of viral genetic, geographic, and host genetic distance matrices.** To

154 determine the relationship between viral genetic, geographic, and host genetic distances, we

155 compiled a data set containing the Kenyan CoV samples generated in this study. The genetic

156 distance matrices were produced from pairwise comparisons either in the form of uncorrected

157 percentage differences or calculated from the phylogenetic trees (patristic distance) using the

158 Patristic v1.0 program (30) The geographic distances (Euclidean distance) were calculated

159 using the formula "distance = (acos((sin(latitude1) * sin(latitude2)) + (cos(latitude1) *

160 cos(latitude2) * cos(longitude2 - longitude1)))) * 6371", with spatial coordinates of the

161 samples derived from the geographic location information.

162      We used Mantel correlation analyses to test the extent of the correlation between

163 these matrices (31). Both simple Mantel's test and partial Mantel's test were performed, and

7

164     the correlation was evaluated with 10000 permutations. To access which of the two factors –

165     geographic or host genetic distance – best explained total variation in the virus genetic

166     distance matrices, we performed multiple linear regression on these distance matrices (32).

167     The statistical significance of each regression was evaluated by performing 10000

168     permutations. To examine whether the degree of virus genetic relatedness corresponded to

169     the scale of geographic distance or host relatedness, we generated Mantel correlograms. In

170     each correlogram, 10-12 distance classes were assigned following an equal-frequency

171     criterion: each class had similar number of pairwise comparisons. All statistical analyses

172     were performed using the Ecodist package implemented in R3.0.2 (33), and all statistical

173     results were considered significant at the $P = 0.05$ level.

174

175     **Full genome sequencing and sequence analyses.** Five viruses representative of the full

176     diversity of the CoVs newly described here were selected for full genome sequencing:

177     BtKY229E-1, BtKY229E-8, BtKYNL63-9a, BtKYNL63-9b, and BtKYNL63-15. We first

178     sequenced a number of conserved regions throughout the genome using several semi-nested

179     or nested consensus degenerate RT-PCR amplicons. These regions were then bridged using

180     sequence-specific RT-PCR followed by Sanger sequencing (< 2 kb) or using the PacBio

181     platform (> 2 kb). The assembled consensus genome sequences from PacBio sequencing

182     were later confirmed by sequence-specific RT-PCR and Sanger sequencing (GenBank

183     accession numbers KY073744-KY073748). The 5' and 3' genome termini were not

184     determined due to the limited RNA remaining, and were derived with PCR primers based on

185     the conserved genome regions in alpha-CoVs.

186         For each complete genome sequence, potential ORFs were predicted based on the

187     conserved core sequence, 5′-CUAAAC-3′, with a minimum length of 66 amino acids.

188     Ribosomal frameshifts were identified based on the presence of the conserved slippery

8

189    sequence, "UUUAAAC". For phylogenetic analyses, the data set was first separated into six

190    ORFs, namely; ORF1a, ORF1b, Spike (S), Envelope (E), Membrane (M), and Nucleoprotein

191    (N) genes. The data set for each gene was translated into amino acid sequences and aligned

192    using MAFFT v7.013. Phylogenetic trees were then inferred using PhyML as described

193    above. Recombination events were first identified from the occurrence of incongruent

194    topologies in these initial phylogenies, and were then confirmed and characterized using

195    Simplot v3.5.1 (34). In the Simplot analysis, seven sequences were analyzed, including the

196    potential recombinant, the parental viruses, as well as an outgroup. The similarity

197    comparisons of recombinant and the other sequences were plotted using a sliding window

198    with a size of 1000 bp and a step size of 10 bp.

199

200    **RESULTS**

201    **Prevalence of CoV in Kenyan bats.** We examined bats from at least 27 species (17 genera)

202    collected over a four year period (2007-2010) from 30 locations across the southern part of

203    Kenya (Figure 1). A total of 2,050 bats samples were screened for CoV RNA using a pan-

204    coronavirus RT-PCR assay. We found an overall prevalence of 11.7% (240/2,050 bats)

205    (Table S1). This overall prevalence is in line with recent reports of CoVs in bats from

206    numerous locations including South Africa, Mexico, Philippines, Kenya, United Kingdom,

207    Japan, Italy, and Ghana (6, 14, 15, 35-40).

208         Bats of the species tested (*Chaerephon pumilus*, *Coleura afra*, *Lissonycteris*

209    *angolensis*, *Miniopterus africanus*, *Neoromicia tenuipinnis*, *Neoromicia* sp*.*, *Nycteris* sp*.*,

210    *Pipistrellus* sp*.*, and *Scotoecus* sp*.*) did not yield CoV positive samples although the sample

211    number was limited and might not reflect the real prevalence (Table S1). Conversely, in bats

212    of several other species the CoV prevalence was high (*Cardioderma cor*, 25%; *Eidolon*

213    *helvum*, 21%; *Epomophorus labiatus*, 28.6%; *Hipposideros* sp*.*, 27.6%; *Miniopterus minor*,

9

214    22.6%; *Otomops martiensseni, 28.6%*; *Rhinolophus hildebrandtii*, 31.3%; *Rhinolophus* sp.,

215    28.9%; *Triaenops afer*, 26.7%). Most species (21/27) were sampled at more than one

216    location. Of note, we detected CoVs in 21% of *E. helvum* bats tested in Kenya, whereas a

217    previous study in Ghana failed to detect any CoVs in a similar number of bats from this

218    species (6).

219

220    **Phylogenetic diversity of Kenyan bat CoVs**. The viral sequences identified in Kenyan bats

221    showed a remarkable diversity within both alpha- and beta-CoVs (Figure 2). Based on our

222    phylogenetic analysis, the CoVs newly identified here can be grouped into 20 phylogenetic

223    lineages (Figure 2). Many of the sampled bat genera are associated with more than one viral

224    lineage. Furthermore, in some cases, the divergence of the CoVs within the same host genera

225    may also be associated with possible differences in sample types. For example, we found two

226    lineages of CoV in *Rousettus aegyptiacus* bats, one of which was present in oral swabs

227    (Figure 2: L7 *Rousettus*) while the other one was identified in fecal swabs (L17 *Rousettus*).

228    The default tissue tropism for bat CoVs is believed to be intestinal and samples of choice are

229    fecal swabs. In agreement with this, only four viruses were identified from oral swab samples

230    (L7 *Rousettus*) as indicated in the phylogeny (Figure 2).

231      Our phylogenetic analyses also revealed a number of cross-species transmission

232    events at the genus level, many of which appeared to be transient spill-overs with no evidence

233    of onward transmission. This pattern was observed as CoV sequences recovered from bats of

234    a particular genus located as tree tips within the phylogenetic diversity that is mainly

235    associated with a different bat genus. From our Kenyan data set, there were seven such cross-

236    species transmission events in total, each represented by a single sequence (dotted red in

237    Figure 2), suggesting these are most likely viruses with limited transmission within new hosts,

238    although this hypothesis requires confirmation on a larger set of samples.

10

239    A more comprehensive and informative phylogeny (Figure 3) was obtained after

240    including the representative global CoV sequences from GenBank, which also included the

241    Kenyan viruses previously reported (15). The phylogeny, which included viral sequences

242    recovered from bats of more than 50 species (30 genera), resulted in an accurate phylogenetic

243    assignment of the viruses described in this study (Figure 3). Importantly, the newly

244    discovered viruses from Kenya have greatly extended our previous work (15) in terms of: (i)

245    expanding the diversity of existing lineages, including the *Miniopterus*, *Rhinolophus*, and

246    *Scotophilus* associated CoV clusters in the genus *Alphacoronavirus*, and the *Rousettus* and

247    *Rhinolophus* associated CoVs clusters in the genus *Betacoronavirus*; and (ii) the discovery of

248    new viruses from either a novel bat host (i.e. *Triaenops*) or new divergent CoV clusters in

249    known hosts (i.e. *Rhinolophus*, *Rousettus*, *Chaerephon*, etc)  (Figure 3).

250    The phylogeny suggests both ancient virus-host co-divergence and recent cross-

251    species transmission of CoVs between bats and other mammalian hosts. The phylogeny

252    clearly demonstrates that CoVs from two host groups, one dominated by bats and the other

253    exclusively by non-chiropteran mammals, formed sister clades for both alpha- and beta-CoVs

254    (Figure 3), suggestive of an ancient divergence between them. Conversely, several non-

255    chiropteran CoVs are nested within the diversity of bat CoVs, suggesting that these viruses

256    are relatively recent introductions from bats. These cross-species transmission events resulted

257    in emergence of severe (SARS-CoV and MERS-CoV) and mild (HCoV-NL63 and HCoV-

258    229E) human pathogens, as well as animal pathogens (Porcine epidemic diarrhea virus

259    [PEDV] and Alpaca respiratory CoV). Interestingly, HCoV-NL63, previously thought to be

260    related to North American tricolored bat (*P. subflavus*) (19), in our phylogeny is deeply

261    nested within the newly identified CoVs from African *Triaenops afer* bats (Figure 3), while

262    the *P. subflavus* virus (labeled green in Figure 3) grouped with a North American CoV

263    sampled from a *Myotis volans* bat (Figure 3). Therefore, *Triaenops afer* bats likely represent

11

264    the most recent chiropteran reservoir host of viruses ancestral to HCoV-NL63. In addition,

265    our results identified 16 additional 229E-like viruses (L14, Figure 2), providing further

266    evidence that *Hipposideros* bats in Africa harbor viruses that are ancestral to HCoV-229E (5,

267    6).

268

269    **Host and spatial dynamics of bat CoVs in Kenya**. We used Mantel's test to compare the

270    virus and host genetic distance matrices, as well as virus and geographic distance matrices.

271    Notably, the correlation values were positive and highly significant in both comparisons

272    (Table 1), suggesting that both host and geography have shaped the structure of virus genetic

273    diversity. This conclusion remained following partial Mantel analyses and multiple linear

274    regression analyses in which we tested the effect between two matrices while controlling for

275    the third (Tables 1 and 2). Importantly, however, in both simple and partial Mantel analyses,

276    the virus genetic distance matrices had much higher correlation with host genetic distance

277    matrices than with geographic distance matrices (Table 1), indicating that bat CoV diversity

278    is more structured by host than by geographic distance.

279         Next, we used Mantel autocorrelograms to examine the effect of (i) geographic

280    distance (Figure 4A) and (ii) host genetic distance on virus diversity (Figure 4B). Host

281    genetic distance decreased from highly significantly positive at short taxonomic distances to

282    highly significantly negative at long distances. Importantly, the crossing-over point was at a

283    host genetic distance of around 0.15-0.19, which marks the boundary of intra- and inter-

284    genera host diversity (Figure 4B). However, no obvious clinal patterns in geographic distance

285    were observed within the Kenyan data set.

286

287    **Full genome characterization and recombination analyses of NL63-like and 229E-like**

288    **viruses**. To further explore the evolution of the NL63-like and 229E-like viruses, we

289    generated the complete genome sequences of five representative bat-derived CoVs: three

290    (BtKYNL63-9a, BtKYNL63-9b, and BtKYNL63-15) were from the NL63-like group and

291    two (BtKY229E-1 and BtKY229E-8) from the 229E-like group (L12-L14, Figure 2). For all

292    the viruses newly described here, the genome structures follow an identical ORF

293    arrangement: ORF1ab-S-ORF4-E-M-N-ORF8 in 229E-related viruses and ORF1ab-S-ORF3-

294    E-M-N-ORFx in NL63-related viruses (Figure 5, Tables 3 and 4). The additional

295    ORF8/ORF$_X$ was identified at the 3' end of the genome in all bat NL63-like and 229E-like

296    viruses characterized in this study, although it was missing in both human viruses (HCoV-

297    229E and HCoV-NL63). The ORF8 in bat 229E-like genomes is named in analogy with the

298    ORF8 of Ghanaian bat and dromedary 229E-like CoVs (5, 12). The ORF8 of BtKY229E-1

299    shared 60% protein identity with its closest relatives while BtKY229E-8 has a shorter and

300    highly divergent ORF8. The ORFx of NL63-like viruses shared very low identity (21-33% at

301    the amino acid level). Similarly to the bat 229E-like CoVs recently discovered in Ghana (5),

302    the S genes in our bat 229E-like CoVs have a considerably longer 5' S1 portion (additional

303    185 amino acids) compared to HCoV-229E and alpaca and dromedary 229E viruses (12).

304            For comparison, we also included 21 genome sequences representative of the

305    diversity in the genus *Alphacoronavirus*. The phylogeny based on the ORF1b protein

306    alignment confirmed that NL63-like and 229E-like groups are monophyletic (Figure 6).

307    Given that each group is associated with a specific bat genus, it is likely that the ORF1b

308    genes of the human viruses (i.e. HCoV-NL63 and HCoV-229E) were ultimately derived from

309    *Triaenops*-associated CoVs and *Hipposideros*-associated CoVs, respectively. The

310    relationship between *Hipposideros* bat CoVs and HCoV-229E was also demonstrated by

311    Corman et al. (5) based on specimens obtained in Ghana. Compared to the viruses described

312    in that study, the newly identified Kenyan viruses (BtKY229E-1 and BtKY229E-8) were

313    among those more distantly related to HCoV-229E (Figure 6 and Table 3). As for the NL63-

13

314    like group, HCoV-NL63 was nested within the diversity of three lineages of *Triaenops*-

315    associated CoVs, among which BtKYNL63-9a showed the closest relationship in all genome

316    regions with the exception of the S gene (Figure 6 and Table 3).

317        Strikingly, the phylogeny of the S protein suggested an entirely different evolutionary

318    history for HCoV-NL63 compared to the rest of the genome (Figure 6). Specifically, for all

319    the proteins with the exception of S, HCoV-NL63 clustered with the NL63-like group.

320    However, in the S protein, HCoV-NL63 was deeply nested within the 229E-like group,

321    associated exclusively with viruses from *Hipposideros* bats, and particularly similar to the

322    sequences BtKY229E-1 and BtKY229E-8 newly identified during this study (Figure 6).

323    Interestingly, BtKY229E-1 exhibited the closest resemblance to HCoV-NL63 in the receptor

324    binding domain (RBD, (41)), especially in the three receptor binding motifs (RBM), whereas

325    other viruses exhibited  less similarity in these regions (Figure 7A). A phylogeny based on

326    the RBD region confirmed our observation (Figure 7B), although it remains uncertain

327    whether these bat viruses utilize the same host cell receptor.

328        To further characterize this recombination event, we performed genome-scale

329    similarity comparisons between HCoV-NL63 and related viruses (Figure 8). The analysis

330    confirmed the chimeric nature of the HCoV-NL63 genome with only the spike protein

331    involved in recombination via two break-points: one located near the 5' end of the S gene and

332    the other at around 200 nucleotides upstream of the 3' end. To exclude the possibility of any

333    artificial recombination, the break-point was further confirmed by RT-PCR and Sanger

334    sequencing, using a single amplicon to cover each break-point. Collectively, these data show

335    that HCoV-NL63 evolved from a recombination event between CoVs from the NL63-like

336    and 229E-like groups.

337        In addition to HCoV-NL63, we identified a number of other recombination events

338    between divergent CoVs involving the S gene. One example is the BtKYNL63-15 newly

339    identified here. Throughout the genome, BtKYNL63-15 showed strong similarity (79% -

340    99% protein identities in the ORF1ab, ORF4, M, E, and N genes) with BtKYNL63-9b. In

341    contrast, the genetic identity between S protein sequences of these viruses was only 53%. In

342    the S protein phylogeny, BtKYNL63-15 did not cluster with NL63-like viruses but instead

343    clustered with *Miniopterus* bat CoV HKU8 and *Chaerophon* bat CoV KY22 (Figure 6).

344    Interestingly, HKU8 itself is a recombinant in the S gene region (Figure 6). These results

345    suggest that the spike protein of CoVs is subject to relatively frequent recombination even

346    between divergent viruses.

347

348    **DISCUSSION**

349       In this study we significantly extended existing knowledge on CoV diversity, their

350    association with specific bat species, the relatedness between bat and human CoVs, and

351    natural recombination events in the CoV spike (S) protein gene between viruses from

352    different lineages.

353       Notably, we found that host species poses a greater influence on CoV diversity in bats

354    than the geographic distance, which can be explained by the ability of bats to fly (including

355    long-distance migrations typical for some species) and disperse their pathogens over vast

356    territories (42). A closer inspection of the Mantel correlogram suggests the presence of less

357    structured (homogenous, mantel statistic r>0), and highly structured (mantel statistic r<0)

358    diversity which, strikingly, corresponds to the division between intra-genera (10% ~ 20%)

359    and inter-genera (> 20%) host genetic distances (Figure 4B). This suggests that within-genus

360    virus transmissions occur significantly more frequently than between-genera transmissions,

361    which is consistent with the previous observations that phylogenetic clustering is less

362    constrained at the host species level than at the genus level (16, 43). While it is commonly

363    accepted that host phylogeny constrains virus cross-species transmission to some extent (44),

15

364    the stronger demarcation at the genus level is of particular interest. In fact, bats of different

365    species, genera, and families frequently roost together (in caves, tree holes, and other

366    shelters), sometimes in dense aggregations, which provide abundant opportunity for

367    mechanical transmission of pathogens between host species. Therefore, our data suggests that

368    distinctions between bats at the genus level might mark a threshold where the differences in

369    cellular and immunological environment become a major challenge for a virus to switch hosts.

370    This, in turn, will lead to the pattern of 'preferential host switching' that has been observed in

371    a number of other viruses (45).

372        The detection of distinctive HCoV-NL63-like and HCoV-229E-like sequences in bats

373    sheds new light on CoV evolution. In particular, we provide strong evidence that HCoV-

374    NL63 has a zoonotic recombinant origin. Although the majority of the HCoV-NL63 genome

375    originates from the viruses circulating in *Triaenops afer* bats, its spike protein gene is derived

376    from a 229E-like virus circulating in *Hipposideros* spp. bats. However, despite the strong

377    signal for recombination, both putative parental strains show substantial genetic distances

378    from human CoVs. This most likely reflects extensive post-recombination sequence

379    divergence, which in turn suggests that the recombination event has occurred prior to the

380    emergence of HCoV-NL63 in humans.

381        Most of the recombination events reported here involve breakpoints around the S

382    gene. Indeed, similar breakpoints are also reported for SARS-CoV and SARS-like CoVs  (24,

383    25), HCoV-OC43 (26), and a feline CoV (46) such that it is seemingly a recombination

384    'hotspot' in many CoVs. It has been argued that a strong secondary structure between ORF1a

385    and S gene may promote transcriptional pulsing, facilitating recombination (47). However,

386    there is also evidence that this recombination hotspot does not exist under non-selective

387    conditions (48), such that it may reflect the successful spread of beneficial recombinants

388     rather than an elevated rate of recombination per se. This hypothesis is supported by the fact

389     that the spike protein is intimately involved in the interaction with the host immune system.

390         Importantly, our results also revealed that recombination has resulted in similar S

391     proteins in the two human viruses HCoV-NL63 and HCoV-229E, such that acquisition of a

392     229E-like S protein may have contributed to the emergence of NL63-like viruses in humans.

393     However, despite this similarity of S protein sequences, these two human viruses utilize

394     different receptors (ACE2 and aminopeptidase-N for HCoV-NL63 and HCoV-229E,

395     respectively) to enter human cells. Within the 229E-like group, the RBD of HCoV-NL63 is

396     more closely related to BtKY229E-8 than to HCoV-229E. The RBD of BtKY229E-8 exhibits

397     greater similarity with that of HCoV-NL63 (Figure 7), and is therefore more likely to be the

398     prototype of RBD in HCoV-NL63.

399         Until recently, most reported recombination events in CoVs involved viruses

400     associated with closely related host species, although recombination between highly

401     divergent CoVs has been demonstrated experimentally (49-51). The apparent lack of

402     interspecies recombination under natural conditions is most likely due to the insufficient

403     collection of complete genome sequences that are truly representative of coronavirus

404     diversity. Indeed, a number of viruses, such as HKU2, display phylogenetic incongruence

405     across different parts of the genome (52), although the lack of one of the putative parental

406     strains has prevented clear identification of a recombinant history.

407         Finally, our study provides insights into the evolution history of CoVs. Although it is

408     unclear whether bats are direct ancestors of all alpha- or beta-CoVs due to the presence of

409     non-bat CoV clades at the basal phylogenetic positions of both genera (Figure 3), bat-borne

410     CoVs constitute a substantial part of the diversities of alpha- or beta-CoVs. In addition, six

411     lineages of non-bat CoVs are nested within the bat-borne clades. These likely represent

412     independent and successful adaptations via shifts from the progenitor reservoir species (bats)

17

413   to other mammals. Four well-characterized human CoVs lie within these clades. However, it

414   is worth noting that bats may not have directly transmitted the viruses to humans. Indeed,

415   HCoV-229E is more closely related to viruses circulating in camels than those in bats,

416   suggesting that camels may be intermediate hosts between bats and humans (12). Similarly,

417   other human CoVs such as SARS-CoV and MERS-CoV all use terrestrial mammals rather

418   than bats as intermediate hosts, which have an increased chance of contact with humans. This

419   underlines a typical zoonotic link of bat-associated CoV to humans via terrestrial mammals.

420

421

422

423   **ACKNOWLEDGEMENTS**

428

429

430

**REFERENCES**

1. **Weiss SR, Leibowitz JL.** 2011. Coronavirus pathogenesis. Adv Virus Res **81:**85-164.

2. **Adams MJ, Lefkowitz EJ, King AM, Bamford DH, Breitbart M, Davison AJ, Ghabrial SA, Gorbalenya AE, Knowles NJ, Krell P, Lavigne R, Prangishvili D, Sanfacon H, Siddell SG, Simmonds P, Carstens EB.** 2015. Ratification vote on taxonomic proposals to the International Committee on Taxonomy of Viruses (2015). Arch Virol **160:**1837-1850.

3. **Lau SK, Woo PC, Li KS, Huang Y, Tsoi HW, Wong BH, Wong SS, Leung SY, Chan KH, Yuen KY.** 2005. Severe acute respiratory syndrome coronavirus-like virus in Chinese horseshoe bats. Proc Natl Acad Sci U S A **102:**14040-14045.

4. **Li W, Shi Z, Yu M, Ren W, Smith C, Epstein JH, Wang H, Crameri G, Hu Z, Zhang H, Zhang J, McEachern J, Field H, Daszak P, Eaton BT, Zhang S, Wang LF.** 2005. Bats are natural reservoirs of SARS-like coronaviruses. Science **310:**676-679.

5. **Corman VM, Baldwin HJ, Tateno AF, Zerbinati RM, Annan A, Owusu M, Nkrumah EE, Maganga GD, Oppong S, Adu-Sarkodie Y, Vallo P, da Silva Filho LV, Leroy EM, Thiel V, van der Hoek L, Poon LL, Tschapka M, Drosten C, Drexler JF.** 2015. Evidence for an Ancestral Association of Human Coronavirus 229E with Bats. J Virol **89:**11858-11870.

6. **Pfefferle S, Oppong S, Drexler JF, Gloza-Rausch F, Ipsen A, Seebens A, Muller MA, Annan A, Vallo P, Adu-Sarkodie Y, Kruppa TF, Drosten C.** 2009. Distant relatives of severe acute respiratory syndrome coronavirus and close relatives of human coronavirus 229E in bats, Ghana. Emerg Infect Dis **15:**1377-1384.

7. **Annan A, Baldwin HJ, Corman VM, Klose SM, Owusu M, Nkrumah EE, Badu EK, Anti P, Agbenyega O, Meyer B, Oppong S, Sarkodie YA, Kalko EK, Lina PH, Godlevska EV, Reusken C, Seebens A, Gloza-Rausch F, Vallo P, Tschapka M, Drosten C, Drexler JF.** 2013. Human betacoronavirus 2c EMC/2012-related viruses in bats, Ghana and Europe. Emerg Infect Dis **19:**456-459.

8. **Lau SK, Li KS, Tsang AK, Lam CS, Ahmed S, Chen H, Chan KH, Woo PC, Yuen KY.** 2013. Genetic characterization of Betacoronavirus lineage C viruses in bats reveals marked sequence divergence in the spike protein of pipistrellus bat coronavirus HKU5 in Japanese pipistrelle: implications for the origin of the novel Middle East respiratory syndrome coronavirus. J Virol **87:**8638-8650.

9. **Corman V, Ithete N, Richards L, Schoeman M, Preiser W, Drosten C, Drexler J.** 2014. Rooting the Phylogenetic Tree of Middle East Respiratory Syndrome Coronavirus by Characterization of a Conspecific Virus from an African Bat. J Virol **88:**11297–11303.

10. **Ithete N, Stoffberg S, Corman V, Cottontail V, Richards L, Schoeman M, Drosten C, Drexler J, Preiser W.** 2013. Close Relative of Human Middle East Respiratory Syndrome Coronavirus in Bat, South Africa. Emerg Infect Dis **19:**1697–1699. .

19

11. **Sabir JS, Lam TT, Ahmed MM, Li L, Shen Y, Abo-Aba SE, Qureshi MI, Abu-Zeid M, Zhang Y, Khiyami MA, Alharbi NS, Hajrah NH, Sabir MJ, Mutwakil MH, Kabli SA, Alsulaimany FA, Obaid AY, Zhou B, Smith DK, Holmes EC, Zhu H, Guan Y.** 2016. Co-circulation of three camel coronavirus species and recombination of MERS-CoVs in Saudi Arabia. Science **351:**81-84.

12. **Corman V, Eckerle I, Memish Z, Liljander A, Dijkman R, Jonsdottir H, Juma NK, Kamau E, Younan M, Al Masri M, Assiri A, Gluecks I, Musa B, Meyer B, Müller M, Hilali M, Bornstein S, Wernery U, Thiel V, Jores J, Drexler J, Drosten C.** 2016. Link of a ubiquitous human coronavirus to dromedary camels. Proc Natl Acad Sci U S A.

13. **O'Shea TJ, Cryan PM, Cunningham AA, Fooks AR, Hayman DT, Luis AD, Peel AJ, Plowright RK, Wood JL.** 2014. Bat flight and zoonotic viruses. Emerg Infect Dis **20:**741-745.

14. **Tao Y, Tang K, Shi M, Conrardy C, Li KS, Lau SK, Anderson LJ, Tong S.** 2012. Genomic characterization of seven distinct bat coronaviruses in Kenya. Virus Res **167:**67-73.

15. **Tong S, Conrardy C, Ruone S, Kuzmin IV, Guo X, Tao Y, Niezgoda M, Haynes L, Agwanda B, Breiman RF, Anderson LJ, Rupprecht CE.** 2009. Detection of novel SARS-like and other coronaviruses in bats from Kenya. Emerg Infect Dis **15:**482-485.

16. **Drexler JF, Corman VM, Drosten C.** 2014. Ecology, evolution and classification of bat coronaviruses in the aftermath of SARS. Antiviral Res **101:**45-56.

17. **Fouchier RA, Hartwig NG, Bestebroer TM, Niemeyer B, de Jong JC, Simon JH, Osterhaus AD.** 2004. A previously undescribed coronavirus associated with respiratory disease in humans. Proc Natl Acad Sci U S A **101:**6212-6216.

18. **Fielding BC.** 2011. Human coronavirus NL63: a clinically important virus? Future Microbiol **6:**153-159.

19. **Huynh J, Li S, Yount B, Smith A, Sturges L, Olsen JC, Nagel J, Johnson JB, Agnihothram S, Gates JE, Frieman MB, Baric RS, Donaldson EF.** 2012. Evidence supporting a zoonotic origin of human coronavirus strain NL63. J Virol **86:**12816-12825.

20. **Corman VM, Rasche A, Diallo TD, Cottontail VM, Stocker A, Souza BF, Correa JI, Carneiro AJ, Franke CR, Nagy M, Metz M, Knornschild M, Kalko EK, Ghanem SJ, Morales KD, Salsamendi E, Spinola M, Herrler G, Voigt CC, Tschapka M, Drosten C, Drexler JF.** 2013. Highly diversified coronaviruses in neotropical bats. J Gen Virol **94:**1984-1994.

21. **Holmes EC.** 2009. The evolution and emergence of RNA viruses, *on* Oxford University Press,.

22. **Lai MM.** 1992. RNA recombination in animal and plant viruses. Microbiol Rev **56:**61-79.

512  23.  **Baric RS, Fu K, Schaad MC, Stohlman SA.** 1990. Establishing a genetic
513      recombination map for murine coronavirus strain A59 complementation groups.
514      Virology **177:**646-656.

515  24.  **Hon CC, Lam TY, Shi ZL, Drummond AJ, Yip CW, Zeng F, Lam PY, Leung FC.**
516      2008. Evidence of the recombinant origin of a bat severe acute respiratory syndrome
517      (SARS)-like coronavirus and its implications on the direct ancestor of SARS
518      coronavirus. J Virol **82:**1819-1826.

519  25.  **Lau SK, Li KS, Huang Y, Shek CT, Tse H, Wang M, Choi GK, Xu H, Lam CS,**
520      **Guo R, Chan KH, Zheng BJ, Woo PC, Yuen KY.** 2010. Ecoepidemiology and
521      complete genome comparison of different strains of severe acute respiratory
522      syndrome-related Rhinolophus bat coronavirus in China reveal bats as a reservoir for
523      acute, self-limiting infection that allows recombination events. J Virol **84:**2808-2819.

524  26.  **Lau SK, Lee P, Tsang AK, Yip CC, Tse H, Lee RA, So LY, Lau YL, Chan KH,**
525      **Woo PC, Yuen KY.** 2011. Molecular epidemiology of human coronavirus OC43
526      reveals evolution of different genotypes over time and recent emergence of a novel
527      genotype due to natural recombination. J Virol **85:**11325-11337.

528  27.  **Pyrc K, Dijkman R, Deng L, Jebbink MF, Ross HA, Berkhout B, van der Hoek L.**
529      2006. Mosaic structure of human coronavirus NL63, one thousand years of evolution.
530      J Mol Biol **364:**964-973.

531  28.  **Katoh K, Standley DM.** 2013. MAFFT multiple sequence alignment software
532      version 7: improvements in performance and usability. Mol Biol Evol **30:**772-780.

533  29.  **Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O.** 2010.
534      New algorithms and methods to estimate maximum-likelihood phylogenies: assessing
535      the performance of PhyML 3.0. Syst Biol **59:**307-321.

536  30.  **Fourment M, Gibbs M.** 2006. PATRISTIC: a program for calculating patristic
537      distances and graphically comparing the components of genetic change. BMC Evol
538      Biol **6:**1-5.

539  31.  **Mantel N.** 1967. The detection of disease clustering and a generalized regression
540      approach. Cancer Res **27:**209-220.

541  32.  **Lichstein JW.** 2007. Multiple regression on distance matrices: a multivariate spatial
542      analysis tool. Plant Ecology **188:**117-131.

543  33.  **Goslee SC, Urban DL.** 2007. The ecodist package for dissimilarity-based analysis of
544      ecological data. Journal of Statistical Software **22:**1-19.

545  34.  **Lole KS, Bollinger RC, Paranjape RS, Gadkari D, Kulkarni SS, Novak NG,**
546      **Ingersoll R, Sheppard HW, Ray SC.** 1999. Full-length human immunodeficiency
547      virus type 1 genomes from subtype C-infected seroconverters in India, with evidence
548      of intersubtype recombination. J Virol **73:**152-160.

549  35.  **Anthony SJ, Ojeda-Flores R, Rico-Chavez O, Navarrete-Macias I, Zambrana-**
550      **Torrelio CM, Rostal MK, Epstein JH, Tipps T, Liang E, Sanchez-Leon M,**
551      **Sotomayor-Bonilla J, Aguirre AA, Avila-Flores R, Medellin RA, Goldstein T,**

21

552        **Suzan G, Daszak P, Lipkin WI.** 2013. Coronaviruses in bats from Mexico. J Gen
553        Virol **94:**1028-1038.

554 36. **August TA, Mathews F, Nunn MA.** 2012. Alphacoronavirus detected in bats in the
555        United Kingdom. Vector Borne Zoonotic Dis **12:**530-533.

556 37. **Balboni A, Palladini A, Bogliani G, Battilani M.** 2011. Detection of a virus related
557        to betacoronaviruses in Italian greater horseshoe bats. Epidemiol Infect **139:**216-219.

558 38. **Geldenhuys M, Weyer J, Nel LH, Markotter W.** 2013. Coronaviruses in South
559        African bats. Vector Borne Zoonotic Dis **13:**516-519.

560 39. **Shirato K, Maeda K, Tsuda S, Suzuki K, Watanabe S, Shimoda H, Ueda N, Iha**
561        **K, Taniguchi S, Kyuwa S, Endoh D, Matsuyama S, Kurane I, Saijo M,**
562        **Morikawa S, Yoshikawa Y, Akashi H, Mizutani T.** 2012. Detection of bat
563        coronaviruses from Miniopterus fuliginosus in Japan. Virus Genes **44:**40-44.

564 40. **Tsuda S, Watanabe S, Masangkay JS, Mizutani T, Alviola P, Ueda N, Iha K,**
565        **Taniguchi S, Fujii H, Kato K, Horimoto T, Kyuwa S, Yoshikawa Y, Akashi H.**
566        2012. Genomic and serological detection of bat coronavirus from bats in the
567        Philippines. Arch Virol **157:**2349-2355.

568 41. **Wu K, Li W, Peng G, Li F.** 2009. Crystal structure of NL63 respiratory coronavirus
569        receptor-binding domain complexed with its human receptor. Proc Natl Acad Sci U S
570        A **106:**19970-19974.

571 42. **Peel AJ, Sargan DR, Baker KS, Hayman DT, Barr JA, Crameri G, Suu-Ire R,**
572        **Broder CC, Lembo T, Wang LF, Fooks AR, Rossiter SJ, Wood JL, AA. C.** 2013.
573        Continent-wide panmixia of an African fruit bat facilitates transmission of potentially
574        zoonotic viruses. Nat Commun **4:**2770.

575 43. **Drexler JF, Gloza-Rausch F, Glende J, Corman VM, Muth D, Goettsche M,**
576        **Seebens A, Niedrig M, Pfefferle S, Yordanov S, Zhelyazkov L, Hermanns U,**
577        **Vallo P, Lukashev A, Muller MA, Deng H, Herrler G, Drosten C.** 2010. Genomic
578        characterization of severe acute respiratory syndrome-related coronavirus in European
579        bats and classification of coronaviruses based on partial RNA-dependent RNA
580        polymerase gene sequences. J Virol **84:**11336-11349.

581 44. **Streicker DG, Turmelle AS, Vonhof MJ, Kuzmin IV, McCracken GF,**
582        **Rupprecht CE.** 2010. Host phylogeny constrains cross-species emergence and
583        establishment of rabies virus in bats. Science **329:**676-679.

584 45. **Charleston MA, Robertson DL.** 2002. Preferential host switching by primate
585        lentiviruses can account for phylogenetic similarity with the primate phylogeny. Syst
586        Biol **51:**528-535.

587 46. **Herrewegh AA, Smeenk I, Horzinek MC, Rottier PJ, de Groot RJ.** 1998. Feline
588        coronavirus type II strains 79-1683 and 79-1146 originate from a double
589        recombination between feline coronavirus type I and canine coronavirus. J Virol
590        **72:**4508-4514.

591 47. **Mills DR, Dobkin C, Kramer FR.** 1978. Template-determined, variable rate of RNA
592 chain elongation. Cell **15:**541-550.

593 48. **Banner LR, Lai MM.** 1991. Random nature of coronavirus RNA recombination in
594 the absence of selection pressure. Virology **185:**441-445.

595 49. **Baric RS, Yount B, Hensley L, Peel SA, Chen W.** 1997. Episodic evolution
596 mediates interspecies transfer of a murine coronavirus. J Virol **71:**1946-1955.

597 50. **Becker MM, Graham RL, Donaldson EF, Rockx B, Sims AC, Sheahan T, Pickles
598 RJ, Corti D, Johnston RE, Baric RS, Denison MR.** 2008. Synthetic recombinant
599 bat SARS-like coronavirus is infectious in cultured cells and in mice. Proc Natl Acad
600 Sci U S A **105:**19944-19949.

601 51. **Masters PS, Rottier PJ.** 2005. Coronavirus reverse genetics by targeted RNA
602 recombination. Curr Top Microbiol Immunol **287:**133-159.

603 52. **Lau SK, Woo PC, Li KS, Huang Y, Wang M, Lam CS, Xu H, Guo R, Chan KH,
604 Zheng BJ, Yuen KY.** 2007. Complete genome sequence of bat coronavirus HKU2
605 from Chinese horseshoe bats revealed a much smaller spike gene with a different
606 evolutionary lineage from the rest of the genome. Virology **367:**428-439.

607

608

609     **Figure Legends**

610     **Figure 1. Map of Kenya showing the geographic locations of 30 bat collection sites.**

611

612     **Figure 2. Phylogeny of RdRp of all CoVs discovered in this study.** The host (bat genus),

613     number of sequences, and operational classification (lineage) are shown on the right of the

614     tree. Branches that represent the minority host genera within the lineage defined by a single

615     dominant host genus are marked with red and labeled with solid circle. The tree is mid-point

616     rooted for clarity only and support values are only shown for internal branches.

617

618     **Figure 3. Phylogenies of RdRp of alphacoronaviruses and betacoronaviruses.** The trees

619     are inferred using representative CoV sequences from this study as well as those obtained

620     from the GenBank. The sequences are labeled with accession number/strain name, host

621     (species) and geographic origin (three letter country code). Different colors are used to

622     distinguish the following groups: Kenyan bat CoVs discovered during this study (orange),

623     CoVs identified from non-bat mammals (blue), the *Perimyotis subflavus* virus previously

624     reported to be related to HCoV-NL63 (green), and the remaining bat viruses (black). The

625     lineage information for Kenyan CoVs is shown to the right of the phylogeny and matches that

626     in Figure 2.

627

628     **Figure 4. Mantel correlograms showing the Kenyan bat CoV RdRp sequences stratified**

629     **by (A) geographic distances and (B) host genetic distances**. A Mantel correlation index (r)

630     was calculated for each of the distance classes. Under the null hypothesis of no relationship

631     between the distance matrices, r values would be close to zero. Positive r values suggest

632     lower genetic distances between case pairs, whereas negative r values suggest higher genetic

633     distances between case pairs. Solid dots: r significantly different from zero; hollow dots: r not

634     significantly different from zero. The graph (B) also shows kernel density plots for intra-

24

635   genus host distances density (grey solid line) and inter-genus host distances density (grey

636   dotted line). The corresponding y-axis for the plot is shown on the right of the figure (B). The

637   grey box in between the two plots marked the transition area between the intra-genus and

638   inter-genus host genetic distances

639

640   **Figure 5. Genome organization of 2 bat 229E-like and 3 bat NL63-like viruses sampled**

641   **from Kenyan bats.** A unified length scale is used for all the genomes. Within each genome,

642   the ORFs (arrow boxes) and ribosomal frame shift sites (vertical lines) are indicated at their

643   corresponding positions.

644

645   **Figure 6. Phylogenetic analyses of major open reading frames of NL63-like and 229E-**

646   **like CoVs in the context of alphacoronaviruses revealing evidence of recombination**.

647   Viruses sequenced in this study are shown in orange. Three potential recombinant genomes

648   of HCoV-NL63, BtKYNL63-15, and HKU8 are indicated with red circles, blue triangles, and

649   brown squares.

650

651   **Figure 7. The relationships between HCoV-NL63 and related viruses at the receptor**

652   **binding domain.** (A) Alignment of NL63-like and 229E-like viruses and related viruses at

653   the receptor binding domain. The positions of three receptor binding motifs (RBMs) are

654   marked with double arrowed line. Residues in the NL63-CoV RBMs that directly contact the

655   ACE2 receptor are marked with red downward arrows. (B) Phylogenetic relationships of

656   NL63-like and 229E-like viruses at receptor binding domain of HCoV-NL63. The tree is

657   based on an amino acid alignment and mid-point rooted.

658

25

659 **Figure 8. Recombination analyses of HCoV-NL63 using Simplot**. Genome-scale similarity

660 comparisons of HCoV-NL63 (query) against BtKYNL63-9a (major parental group, blue),

661 BtKYNL63-9b (green), BtKY229E-8 (minor parental group, red), HCoV-229E (orange),

662 BtCoV/FO1A-F2/Hip_aba/GHA/2010 (pink), and Alaca respiratory CoV (brown). A full

663 genome structure, with reference to HCoV-NL63, is shown above the similarity plot, marking

664 the positions and boundaries of the major open reading frames. At the beginning of the S

665 gene, the flat-line followed by a sudden drop of similarity is due to a gap (deletion within

666 HCoV-229E S gene) in the alignment.

667

668
669
670
671 **Tables**

672

673 **Table 1.** Results of Mantel tests and partial Mantel tests comparing two factors (host genetic
674 distance and geographic distance) that predict the structure of virus genetic diversity

| Model | $r$ value for Kenyan bats ($P$ value) |
|---|---|
| Host[a] | 0.5265 ($P < 0.0001$)[c] |
| Host \| Geography[b] | 0.5055 ($P < 0.0001$)[c] |
| Geography[a] | 0.2122 ($P < 0.0001$)[c] |
| Geography \| Host[b] | 0.1285 ($P = 0.0005$)[c] |

675 [a]Mantel test; [b]partial Mantel test; [c]significant at 0.001.

676

677

678 **Table 2.** Multiple regression of virus genetic distance against host genetic distance and
679 geographic distance in Kenyan bat CoVs (2007-2010)

| Variable | Correlation coefficient | $P$-value |
|---|---|---|
| Host | 7.58E-01 | 1.00E-04 |
| Geography | 1.19E-06 | 1.00E-02 |

680

26

681 **Table 3. Sequence comparisons of the Kenyan bat CoVs with HCoV-229E or HCoV-NL63**

| | Genome identity | Concatenated domains | ADRP | nsp5 | nsp12 | nsp13 | nsp14 | nsp15 | nsp16 | 1ab | S | ORF3/4 | E | M | N |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Nucleotide % | | | | | | Amino acid % identity to HCoV-229E | | | | | | | | |
| **BtKY229E-1** | 88 | 98 | 92 | 98 | 97 | 99 | 97 | 96 | 94 | 95 | 75 | 92 | 93 | 90 | 78 |
| **BtKY229E-8** | 88 | 97 | 89 | 98 | 98 | 98 | 97 | 97 | 94 | 96 | 74 | 94 | 97 | 90 | 68 |
| | Nucleotide % | | | | | | Amino acid % identity to HCoV-NL63 | | | | | | | | |
| **BtKYNL63-9a** | 78 | 91 | 75 | 89 | 93 | 94 | 89 | 88 | 94 | 86 | 53 | 67 | 80 | 82 | 69 |
| **BtKYNL63-9b** | 68 | 83 | 51 | 76 | 88 | 91 | 82 | 81 | 84 | 72 | 52 | 55 | 64 | 61 | 51 |
| **BtKYNL63-15** | 68 | 84 | 51 | 76 | 88 | 91 | 82 | 81 | 87 | 72 | 49 | 55 | 62 | 58 | 52 |

682

683

684

685

686

27

687    **Table 4. Genomic features of the open reading frames (ORF) in the Kenyan bat CoVs and their putative transcription regulatory sequences (TRS).**

| Virus | | 229E-like virus | | NL63-like virus | | |
|---|---|---|---|---|---|---|
| | | BtKY229E-1 | BtKY229E-8 | BtKYNL63-9a | BtKYNL63-9b | BtKYNL63-15 |
| Sequences (nt) | | 27837 | 27666 | 28363 | 28677 | 28479 |
| GC% | | 38 | 39 | 39 | 43 | 43 |
| ORF1ab (nt) | ORF size (nt) | 20286 | 20304 | 20277 | 20349 | 20355 |
| | Putative TRS | TCTCAACTAAAC(N219)AUG | TCTCAACTAAAC(N219)AUG | TCAACTAAAC(N214)AUG | CTCAACTAAAC(N215)AUG | TCTCAACTAAAC(N215)AUG |
| S | ORF size (nt) | 4095 | 4095 | 4119 | 4122 | 4134 |
| | Putative TRS | TCTCAACTAAATAAAAUG | UCTCAACUAAA(4)AUG | TCAACTAAAC(N1)AUG | CTCAACTAAAUG | TCAACTAAAC(N1)AUG |
| ORF3/4 | ORF size (nt) | 681 | 684 | 684 | 684 | 684 |
| | Putative TRS | TCAACTAAAC(N37)AUG | TCAACTAAAC(N37)AUG | TCAACTAAAC(N37)AUG | TCAACTAAAC(N37)AUG | CAACUAAAC(N37)AUG |
| E | ORF size (nt) | 234 | 234 | 234 | 234 | 234 |
| | Putative TRS | TCTCAACTAAAC(N151)AUG | TCTTCAATGTAAC(N281)AUG | TTATAAC(N79)AUG | TCTGCTAAC(N151)AUG | TCTGATAAC(N151)AUG |
| M | ORF size (nt) | 681 | 681 | 693 | 681 | 684 |
| | Putative TRS | CTAAACTAAAC(N4)AUG | CTAAACTAAAC(N4)AUG | CTAAAC(N6)AUG | TCTAAACTAAAC(N4)AUG | UCUAAACUAAA(N4)AUG |
| N | ORF size (nt) | 1161 | 1200 | 1225 | 1302 | 1302 |
| | Putative TRS | TTAATCTAAAC(N11)AUG | ATCTAAAC(N11)AUG | TCTAAACTAAAC(N3)AUG | CTAAACCAAAC(N4)AUG | UCUAAACUAAAC(N4)AUG |
| ORF8/ORFx | ORF size (nt) | 288 | 198 | 287 | 291 | 270 |
| | Putative TRS | UCAACUAAAAC(1)AUG | UCAACUAAAAC(4)AUG | CAAAACCUAAC(N12)AUG | TCAACTAAAC(N567)AUG | CAACUAAAC(N234)AUG |

688

Ethiopia

Uganda

N

1, 2

27, 36-38

5

19

Equator

18  4

45

20

7

40, 41

42

25

10

23

33

24

Tanzania

26

22

43

35

16, 21

13, 14

44

*Rousettus*

*Rousettus*

*Rhinolophus*

L1 *Minopterus* (n=55)

90

100

91

L2 *Minopterus* (n=36)

*Hipposideros*

82

99 L3 *Rhinolophus* (n=8)

L4 *Cardioderma* (n=2)

88 L5 *Rhinolophus* (n=3)

85 *Rousettus*

L6 *Rhinolophus* (n=11)

L7 *Rousettus* (Oral samples, n=4)

80 *Epomophorus*

L8 *Otomops* (n=5)

91 95

L9 *Otomops* (n=6)

98

L10 *Chaerephon* (n=12)

*Epomophorus*

91

L11 *Scotophilus* (n=1)

86 99 L12 *Triaenops* (n=8)

L13 *Triaenops* (n=1)

81 100

L14 *Hipposideros* (n=16)

*Alpha-CoV*

L15 *Chaerephon* (n=1)

*Beta-CoV*

98 L16 *Rhinolophus* (n=3)

100 L17 *Rousettus* (n=12)

L18 *Epomophorus* (n=11)

100

100

95 L19 *Eidolon* (n=38)

0.5

L20 *Epomophorus* (n=1)

## Alpha-CoV

```
95   HQ728484 KEN Miniopterus sp
99     BtKY195 Miniopterus minor              L1
       BtKY130 Miniopterus minor
92     GU190240 BGR Miniopterus schreibersii
       DQ666337 CHN Miniopterus magnater
99     BtKY258 Miniopterus minor              L2
     95  BtKY224 Miniopterus minor
         DQ249228 CHN Miniopterus pusillus
87       DQ666339 CHN Miniopterus magnater
         EU834956 AUT Miniopterus australis
99       KF515987 NZL Mystacina tuberculata
         EU834951 AUT Myotis macropus
91 91    HQ184058 ESP Pipistrellus kuhlii
         KJ473809 CHN Nyctalus velutinus
         KF843855 ZAF Neoromicia cf capensis
95       GU190239 BGR Nyctalus leisleri
82       KT345294 FRA Pipistrellus pipistrellus
         JQ731775 CRI Anoura geoffroyi
         JQ731784 PAN Artibeus jamaicensis
90     99  HQ728480 KEN Cardioderma cor
           BtKY242 Cardioderma cor              L4
           BtKY236 Rhinolophus landeri          L3
        99  100  GU190233 BGR Rhinolophus fer.
        85       BtKY244 Rhinolophus hilderbrandtii   L5-6
                 BtKY70 Rhinolophus sp
                 KU343197 CHN Rhinolophus affinis
                 DQ648854 CHN Rhinolophus sp
             99    BtKY117 Rousettus aegyptiacus  L7
                   JQ989272 CHN Hipposideros sp
                   JQ989270 CHN Rousettus sp
        91   DQ648823 CHN Scotophilus kuhlii
        87     BtKY280 Scotophilus dingani       L11
               KF569988 CHN Myotis davidii
               KF294382 CHN Myotis davidii
        87   JF440355 GBR Myotis nattereri
             JF440350 GBR Myotis daubentonii
        94     EU375868 DEU Pipistrellus pygmaeus
        97     EU375864 DEU Pipistrellus nathusii
               HM368166 DEU Myotis myotis
        89     DQ249224 CHN Myotis ricketti
               EF544565 USA Myotis occultus
               EF185992 Porcine epidemic diarrhea virus
        90     KU182966 CHN Murina leucogaster
        97     KF294376 CHN Murina leucogaster
             97   EF544566 USA Eptesicus fuscus
        98        JX537914 USA Eptesicus fuscus
                  JQ731799 BRA Molossus rufus
                  KF569991 CHN Myotis davidii
        94        KJ473806 CHN Myotis ricketti
                  HQ184050 ESP Myotis blythii
        90   93   HQ336976 USA Myotis volans
                  JX537913 USA Perimyotis subflavus
                  KC110771 BRA NA
        92     BtKY273 Otomops martiensseni      L8
               BtKY147 Chaerophon sp
        98     HQ728486 KEN Chaerophon sp
             98   BtKY275 Otomops martiensseni    L9-10
             88   BtKY270 Chaerophon sp
             82   BtKY204 Epomophorus labiatus
             85   JQ410000 Alpaca respiratory coronavirus
                  BtKY229E-1 Hipposideros sp
                  NC002645 Human coronavirus 229E
        94        BtKY229E-8 Hipposideros vittatus   L14
                  JX174639 GAB Hipposideros caffer
                  KT253270 GHA Hipposideros abae
        100       FJ710045 GHA Hipposideros sp
                  FJ710044 GHA Hipposideros sp
        92     97   NC005831 Human coronavirus NL63
                    BtKYNL63-9a Triaenops afer
                    BtKYNL63-15 Triaenops afer       L12-13
                    BtKYNL63-9b Triaenops afer
        99   99   HQ728481 KEN Chaerophon sp
                  BtKY210 Chaerephon sp              L15
                  JQ731790 CRI Carollia perspicillata
                  JQ731788 PAN Artibeus lituratus
        96        JQ731782 PAN Phyllostomus discolor
                  EU769558 TTO Glossophaga soricina
93
96   Alphacoronavirus1, include:
     Feline coronavirus, Transmissible gastroenteritis virus, Rodent CoVs etc
```

0.5

## Beta-CoV

```
          GQ153539 CHN Rhinolophus sinicus
97        DQ412042 CHN Rhinolophus ferrumequinum
          DQ412043 CHN Rhinolophus macrotis
          AY304488 SARS Coronavirus SZ16
84 89     DQ071615 CHN Rhinolophus pearsoni
100       BtKY237 Rhinolophus hilderbrandtii    L16
94        GU190227 BGR Rhinolophus mehelyi
       95  GU190231 BGR Rhinolophus ferrumequinum
          GQ404795 SVN Rhinolophus hipposideros
          EU834950 AUT Rhinonycteris aurantius
86        KU343200 CHN Hipposideros pomona
          KF636752 CHN Hipposideros pratti
          HQ166910 NGA Hipposideros commersoni
97        JX174638 GAB Hipposideros caffer
94        FJ710050 GHA Hipposideros sp
          FJ710043 GHA Hipposideros sp
          JX899384 GHA Nycteris sp
99        KC545386 DEU Erinaceus europaeus
          KC545383 DEU Erinaceus europaeus
          KC243390 ROU Pipistrellus pygmaeus
99        KC243392 UKR Pipistrellus nathusii
          KU740200 MERS CoV/camel/Egypt/NRCE-NC163/2014
          JX869059 MERS Coronavirus EMC2012
90   100   DQ648794 CHN Tylonycteris pachypus
85        KU182965 CHN Myotis daubentonii
          DQ648819 CHN Pipistrellus pipistrellus
          EF065509 CHN Pipistrellus abramus
97        DQ249221 CHN Pipistrellus abramus
88        DQ648809 CHN Pipistrellus abramus
          KC633197 CRI Carollia perspicillata
90        KC886322 MEX Pteronotus davyi
98        KC633195 CRI Pteronotus parnellii
          BtKY92 Eidolon helvum
94        HQ728482 KEN Eidonlon sp                L19
          BtKY89 Eidolon helvum
89        KU131211 NGA Eidolon helvum
          BtKY54 Epomophorus labiatus
          BtKY182 Epomophorus labiatus            L18 & 20
          BtKY55 Epomophorus labiatus
99        BtKY234 Epomophorus labiatus
98   89     AB539081 PHL Cynopterus brachyotis
     81     KU182962 CHN Cynopterus sphinx
            AB683970 PHL Ptenochirus jagori
            AB918719 IDN Dobsonia moluccensis
     92     HM211100 CHN Rousettus leschenaulti
            BtKY76 Rousettus aegyptiacus
     90     BtKY221 Rousettus aegyptiacus          L17
     96     HQ728483 KEN Rousettus sp
     87     EF065513 CHN Rousettus leschenaultii
     93     EF065516 CHN Rousettus leschenaultii
100   Betacoronavirus lineage A, include:
      Murine hepatitis virus, Bovine coronavirus,
      Rabbit coronavirus, OC43 etc
```
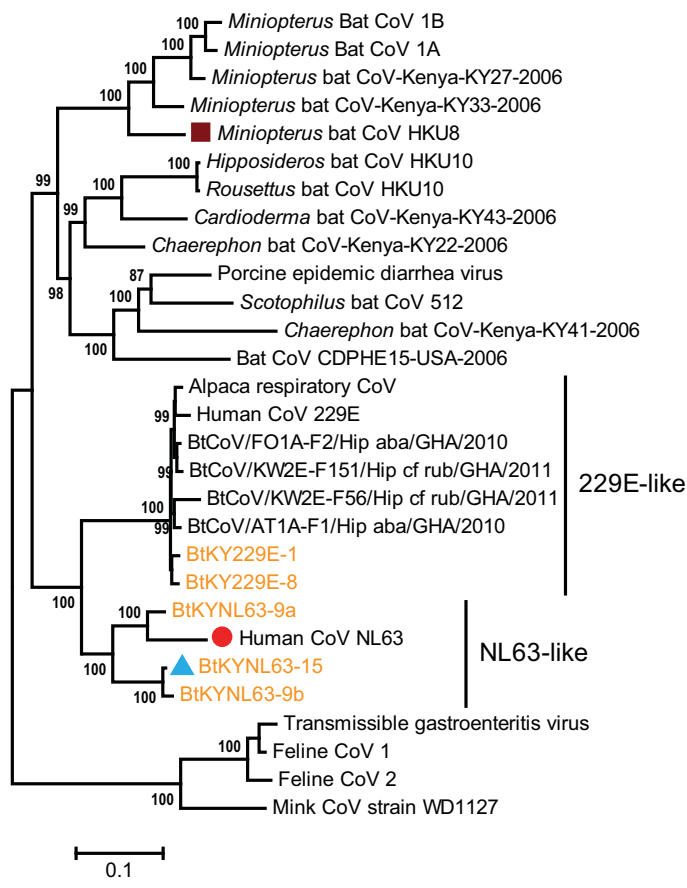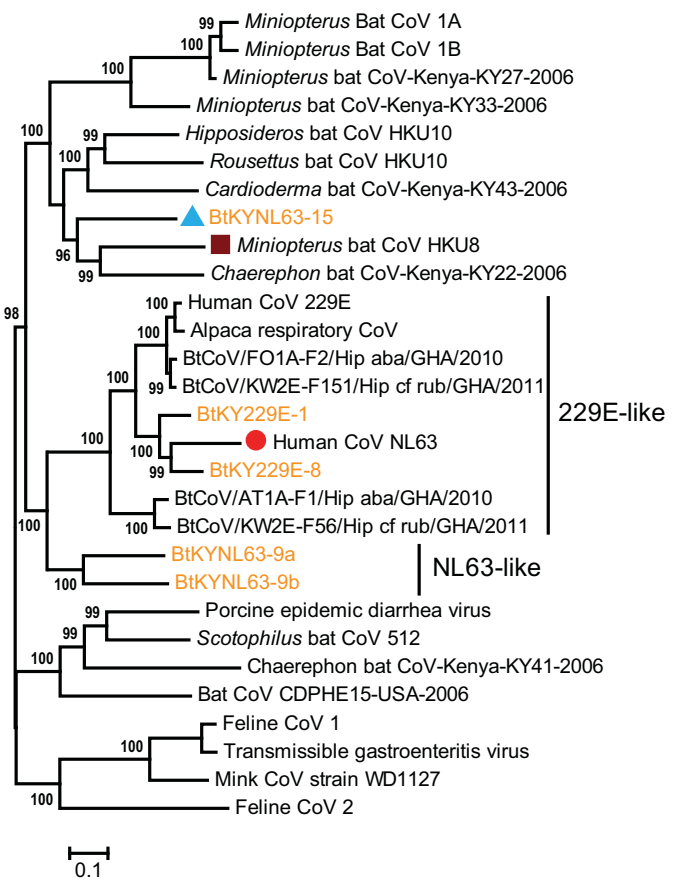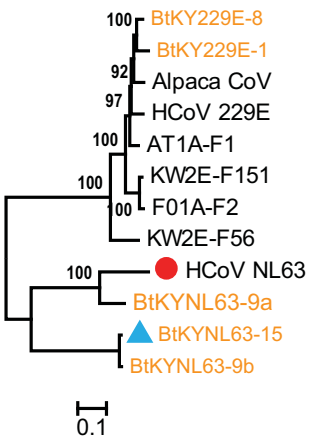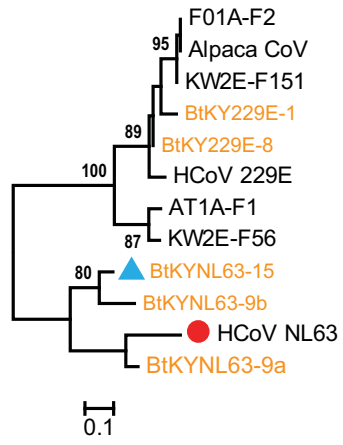
0.5

## 1b Protein



## Spike Protein (non-recombinant region)



**1a**



**E**



**M**



**N**

ORF1ab  S  N

Window : 1000 bp, Step: 10 bp, GapStrip: On, Kimura (2-parameter), T/t: 2.0

HCoV-NL63 vs BtKYNL63-9a
HCoV-NL63 vs BtKYNL63-9b
HCoV-NL63 vs BtKY229E-8
HCoV-NL63 vs Human coronavirus 229E
HCoV-NL63 vs BtCoV/FO1A-F2/Hip aba/GHA/2010
HCoV-NL63 vs Alpaca respiratory coronavirus

Genetic Similarity

Position (bp)