

## 1 **Origin and cross-species transmission of bat coronaviruses in China**

2 Alice Latinne<sup>1§¶</sup>, Ben Hu<sup>2¶</sup>, Kevin J. Olival<sup>1</sup>, Guangjian Zhu<sup>1</sup>, Libiao Zhang<sup>3</sup>, Hongying Li<sup>1</sup>, Aleksei A.

3 Chmura<sup>1</sup>, Hume E. Field<sup>1,4</sup>, Carlos Zambrana-Torrel<sup>1</sup>, Jonathan H. Epstein<sup>1</sup>, Bei Li<sup>2</sup>, Wei Zhang<sup>2</sup>, Lin-Fa

4 Wang<sup>5</sup>, Zheng-Li Shi<sup>2\*</sup>, Peter Daszak<sup>1\*</sup>

5 <sup>1</sup>EcoHealth Alliance, New York, USA;

6 <sup>2</sup>Key laboratory of special pathogens and biosafety, Wuhan Institute of Virology, Center for Biosafety  
7 Mega-Science, Chinese Academy of Sciences, Wuhan, China;

8 <sup>3</sup>Guangdong Institute of Applied Biological Resources, Guangdong Academy of Sciences, Guangzhou,  
9 China;

10 <sup>4</sup>School of Veterinary Science, The University of Queensland, Brisbane, Australia.

11 <sup>5</sup>Programme in Emerging Infectious Diseases, Duke-NUS Medical School, Singapore.

12

13 <sup>§</sup>Current Address: Wildlife Conservation Society, Viet Nam Country Program, Ha Noi, Viet Nam; Wildlife  
14 Conservation Society, Health Program, Bronx, New York, USA;

15 <sup>¶</sup>Authors contributed equally to this paper

16 \*Correspondence should be addressed to: [daszak@ecohealthalliance.org](mailto:daszak@ecohealthalliance.org); [zlishi@wh.iov.cn](mailto:zlishi@wh.iov.cn)

17

## 18 **Abstract**

19 Bats are presumed reservoirs of diverse coronaviruses (CoVs) including progenitors of Severe Acute  
20 Respiratory Syndrome (SARS)-CoV and SARS-CoV-2, the causative agent of COVID-19. However, the  
21 evolution and diversification of these coronaviruses remains poorly understood. We used a Bayesian  
22 statistical framework and sequence data from all known bat-CoVs (including 630 novel CoV sequences)

23 to study their macroevolution, cross-species transmission, and dispersal in China. We find that host-  
24 switching was more frequent and across more distantly related host taxa in alpha- than beta-CoVs, and  
25 more highly constrained by phylogenetic distance for beta-CoVs. We show that inter-family and -genus  
26 switching is most common in Rhinolophidae and the genus *Rhinolophus*. Our analyses identify the host  
27 taxa and geographic regions that define hotspots of CoV evolutionary diversity in China that could help  
28 target bat-CoV discovery for proactive zoonotic disease surveillance. Finally, we present a phylogenetic  
29 analysis suggesting a likely origin for SARS-CoV-2 in *Rhinolophus* spp. bats.

30

## 31 Introduction

32 Coronaviruses (CoVs) are RNA viruses causing respiratory and enteric diseases with varying  
33 pathogenicity in humans and animals. All CoVs known to infect humans are zoonotic, or of animal origin,  
34 with many thought to originate in bat hosts<sup>1,2</sup>. Due to their large genome size (the largest non-  
35 segmented RNA viral genome), frequent recombination and high genomic plasticity, CoVs are prone to  
36 cross-species transmission and are able to rapidly adapt to new hosts<sup>1,3</sup>. This phenomenon is thought to  
37 have led to the emergence of a number of CoVs affecting livestock and human health<sup>4-9</sup>. Three of these  
38 causing significant outbreaks originated in China during the last two decades. Severe Acute Respiratory  
39 Syndrome (SARS)-CoV emerged first in humans in Guangdong province, southern China, in 2002 and  
40 spread globally, causing fatal respiratory infections in close to 800 people<sup>10-12</sup>. Subsequent investigations  
41 identified horseshoe bats (genus *Rhinolophus*) as the natural reservoirs of SARS-related CoVs and the  
42 likely origin of SARS-CoV<sup>13-16</sup>. In 2016, Swine Acute Diarrhea Syndrome (SADS)-CoV caused the death of  
43 over 25,000 pigs in farms within Guangdong province<sup>17</sup>. This virus appears to have originated within  
44 *Rhinolophus* spp. bats, and belongs to the HKU2-CoV clade previously detected in bats in the region<sup>17-19</sup>.  
45 In 2019, a novel coronavirus (SARS-CoV-2) caused an outbreak of respiratory illness (COVID-19) first  
46 detected in Wuhan, Hubei province, China, which has since become a pandemic. This emerging human  
47 virus is closely related to SARS-CoV, and also appears to have originated in horseshoe bats<sup>20,21</sup> - with its  
48 full genome 96% similar to a viral sequence reported from *Rhinolophus affinis*<sup>20</sup>. Closely related  
49 sequences were also identified in Malayan pangolins<sup>22,23</sup>.

50 A growing body of research has identified bats as the evolutionary sources of SARS- and Middle East  
51 Respiratory Syndrome (MERS)-CoVs<sup>13,14,24-26</sup>, and as the source of progenitors for the human CoVs, NL63  
52 and 229E<sup>27,28</sup>. The emergence of SARS-CoV-2 further underscores the importance of bat-origin CoVs to  
53 global health, and understanding their origin and cross-species transmission is a high priority for

54 pandemic preparedness<sup>20,29</sup>. Bats harbor the largest diversity of CoVs among mammals and two CoV  
55 genera, alpha- and beta-CoVs ( $\alpha$ - and  $\beta$ -CoVs), have been widely detected in bats from most regions of  
56 the world<sup>30,31</sup>. Bat-CoV diversity seems to be correlated with host taxonomic diversity globally, the  
57 highest CoV diversity being found in areas with the highest bat species richness<sup>32</sup>. Host switching of  
58 viruses over evolutionary time is an important mechanism driving the evolution of bat coronaviruses in  
59 nature and appears to vary geographically<sup>32,33</sup>. However, detailed analyses of host-switching have been  
60 hampered by incomplete or opportunistic sampling, typically with relatively low numbers of viral  
61 sequences from any given region<sup>34</sup>.

62 China has a rich bat fauna, with more than 100 described bat species and several endemic species  
63 representing both the Palearctic and Indo-Malay regions<sup>35</sup>. Its situation at the crossroads of two  
64 zoogeographic regions heightens China's potential to harbor a unique and distinctive CoV diversity.  
65 Since the emergence of SARS-CoV in 2002, China has been the focus of an intense viral surveillance and  
66 a large number of diverse bat-CoVs has been discovered in the region<sup>36-44</sup>. However, the macroevolution  
67 of CoVs in their bat hosts in China and their cross-species transmission dynamics remain poorly  
68 understood.

69 In this study, we analyze an extensive field-collected dataset of bat-CoV sequences from across China.  
70 We use a phylogeographic Bayesian statistical framework to reconstruct virus transmission history  
71 between different bat host species and virus spatial spread over evolutionary time. Our objectives were  
72 to compare the macroevolutionary patterns of  $\alpha$ - and  $\beta$ -CoVs and identify the hosts and geographical  
73 regions that act as centers of evolutionary diversification for bat-CoVs in China. These analyses aim to  
74 improve our understanding of how CoVs evolve, diversify, circulate among, and transmit between bat  
75 families and genera to help identify bat hosts and regions where the risk of CoV spillover is the highest.

## 76 **Results**

## 77 **Taxonomic and geographic sampling**

78 We generated 630 partial sequences (440 nt) of the *RNA-dependent RNA polymerase (RdRp)* gene from  
79 bat rectal swabs collected in China and added 608 bat-CoV and eight pangolin CoV sequences from  
80 China available in GenBank or GISAID to our datasets (list of GenBank and GISAID accession numbers  
81 available in Supplementary Note 1). For each CoV genus, two datasets were created: one including all  
82 bat-CoV sequences with known host (host dataset) and one including all bat-CoV sequences with known  
83 sampling location at the province level (geographic dataset). To create a geographically discrete  
84 partitioning scheme that was more ecologically relevant than administrative borders for our  
85 phylogeographic reconstructions, we defined six zoogeographic regions within China by clustering  
86 provinces with similar mammalian diversity using hierarchical clustering<sup>45</sup> (see Methods): South western  
87 region (SW), Northern region (NO), Central northern region (CN), Central region (CE), Southern region  
88 (SO) and Hainan island (HI) (Fig. 1 and Supplementary Fig. 1).

89 Our host datasets included 701  $\alpha$ -CoV sequences (353 new sequences, including 102 new SADSr-CoV  
90 sequences (*Rhinacovirus*) from 41 bat species (14 genera, five families) and 528  $\beta$ -CoV sequences (273  
91 new sequences, including 97 new SARSr-CoV sequences (*Sarbecovirus*) from 31 bat species (15 genera,  
92 four families) (Supplementary Table 1). Our geographic datasets included 677  $\alpha$ -CoV sequences from six  
93 zoogeographic regions (22 provinces) and 503  $\beta$ -CoV sequences from five zoogeographic regions (21  
94 provinces) (Fig. 1). As some regions or hosts were overrepresented in our datasets, we also created and  
95 ran our analyses using a more uniform subset of our sequence data that included ~30 randomly-selected  
96 sequences per host family or region to mitigate sampling and surveillance intensity bias.

## 97 **Ancestral hosts and cross-species transmission**

98 We used a Bayesian discrete phylogeographic approach implemented in BEAST<sup>46</sup> to reconstruct the  
99 ancestral host of each node in the phylogenetic tree using bat host family as a discrete character state.

100 The phylogenetic reconstructions for  $\alpha$ -CoVs in China suggest an evolutionary origin within rhinolophid  
101 and vespertilionid bats (Fig. 2A). The first  $\alpha$ -CoV lineage to diverge historically corresponds to the  
102 subgenus *Rhinacovirus* (L1), originating within rhinolophid bats, and includes sequences related to  
103 HKU2-CoV and SADS-CoV (Supplementary Fig. 2). Then several lineages, labelled L2 to L7, emerged from  
104 vespertilionid bats (Fig. 2A). The subgenus *Decacovirus* (L2) includes sequences mostly associated with  
105 the Rhinolophidae and Hipposideridae and related to HKU10-CoV (Supplementary Fig. 3), while the  
106 subgenera *Myotacovirus* (L3) and *Pedacovirus* (L5) as well as an unidentified lineage (L4) include CoVs  
107 mainly from vespertilionid bats and related to HKU6-, HKU10-, and 512-CoVs (Supplementary Fig. 4-5).  
108 Finally, a well-supported node comprises the subgenera *Nyctacovirus* (L6) from vespertilionid bats and  
109 *Minunacovirus* (L7) from miniopterid bats, and includes HKU7-, HKU8-, 1A-, and 1B-CoVs  
110 (Supplementary Fig. 6). These seven  $\alpha$ -CoV lineages are mostly associated with a single host family but  
111 each also included several sequences identified from other bat families (Fig. 2A, Supplementary Fig. 2-6  
112 and Supplementary Table 1), suggesting frequent cross-species transmission events have occurred  
113 among bats. Ancestral host reconstructions based on the random data subset, to normalize sampling  
114 effort, gave very similar results with rhinolophids and vespertilionids being the most likely ancestral  
115 hosts of most  $\alpha$ -CoV lineages too (Supplementary Fig. 7A). However, the topology of the tree based on  
116 the random subset was slightly different as the lineage L5 was paraphyletic.

117 Chinese  $\beta$ -CoVs likely originated from vespertilionid and rhinolophid bats (Fig. 2B). The MCC tree was  
118 clearly structured into four main lineages: *Merbecovirus* (Lineage C), including MERS-related (MERSr-)  
119 CoVs, HKU4- and HKU5-CoVs and strictly restricted to vespertilionid bats (Supplementary Fig. 8);  
120 *Nobecovirus* (lineage D), originating from pteropodid bats and corresponding to HKU9-CoV  
121 (Supplementary Fig. 9); *Hibecovirus* (lineage E) comprising sequences isolated in hipposiderid bats  
122 (Supplementary Fig. 10) and *Sarbecovirus* (Lineage B) including sequences related to HKU3- and SARS-  
123 related (SARSr-) CoVs originating in rhinolophid bats (Supplementary Fig. 11). We show that SARS-CoV-2

124 forms a divergent clade within *Sarbecovirus* and is most closely related to viruses sampled from  
125 *Rhinolophus malayanus* and *R. affinis* and from Malayan pangolins (*Manis javanica*) (Fig. 3). Similar tree  
126 topology and ancestral host inference were obtained with the random subset (Supplementary Fig. 7B).  
127 We used a Bayesian Stochastic Search Variable Selection (BSSVS) procedure<sup>47</sup> to identify viral host  
128 switches (transmission over evolutionary time) between bat families and genera that occurred along the  
129 branches of the MCC annotated tree and calculated Bayesian Factor (BF) to estimate the significance of  
130 these switches (Fig. 4). We identified nine highly supported (BF > 10) inter-family host switches for  $\alpha$ -  
131 CoVs and three for  $\beta$ -CoVs (Fig. 4A and 4B). These results are robust over a range of sample sizes, with  
132 seven of these nine switches for  $\alpha$ -CoVs and the exact same three host switches for  $\beta$ -CoVs having  
133 strong BF support (BF > 10) when analyzing our random subset (Supplementary Tables 2 and 3). To  
134 quantify the magnitude of these host switches, we estimated the number of host switching events  
135 (Markov jumps)<sup>48,49</sup> along the significant inter-family switches (Fig. 4C and 4D) and estimated the rate of  
136 inter-family host switching events per unit of time for each CoV genus. The rate of inter-family host  
137 switching events was five times higher in the evolutionary history of  $\alpha$ - (0.010 host switches/unit time)  
138 than  $\beta$ -CoVs (0.002 host switches/unit time) in China. For  $\alpha$ -CoVs, host switching events from the  
139 Rhinolophidae and the Miniopteridae were greater than from other bat families while rhinolophids were  
140 the highest donor family for  $\beta$ -CoVs. The Rhinolophidae and the Vespertilionidae for  $\alpha$ -CoVs and the  
141 Hipposideridae for  $\beta$ -CoVs received the highest numbers of switching events (Fig. 4C and 4D). When  
142 using the random dataset, similar results were obtained for  $\beta$ -CoVs while rhinolophids were the highest  
143 donor family for  $\alpha$ -CoVs (Supplementary Tables 4 and 5).

144 At the genus level, we identified 20 highly supported inter-genus host switches for  $\alpha$ -CoVs, 17 of them  
145 were also highly significant using the random subset (Fig. 5A and Supplementary Table 6). Sixteen highly  
146 supported inter-genus switches were identified for  $\beta$ -CoVs (Fig. 5B). Similar results were obtained for  
147 the random  $\beta$ -CoV subset (Supplementary Table 7). Most of the significant cross-genus CoV switches for

148  $\alpha$ -CoVs, 15 of 20 (75%), were between genera in different bat families, while this proportion was only 6  
149 of 16 (37.5%) for  $\beta$ -CoVs. The estimated rate of inter-genus host switching events (Markov jumps) was  
150 similar for  $\alpha$ - (0.014 host switches/unit time) and  $\beta$ -CoVs (0.014 host switches/unit time). For  $\alpha$ -CoVs,  
151 *Rhinolophus* and *Miniopterus* were the greatest donor genera and *Rhinolophus* was the greatest receiver  
152 (Supplementary Table 8). For  $\beta$ -CoVs, *Rousettus* was the greatest donor and *Eonycteris* the greatest  
153 receiver genus (Supplementary Table 9).

#### 154 **CoV spatiotemporal dispersal in China**

155 We used our Bayesian discrete phylogeographic model with zoogeographic regions as character states  
156 to reconstruct the spatiotemporal dynamics of CoV dispersal in China. Eleven and seven highly  
157 significant (BF > 10) dispersal routes within China were identified for  $\alpha$ - and  $\beta$ -CoVs, respectively (Fig. 6).  
158 Seven and five of these dispersal routes, respectively, remained significant when using our random  
159 subsets (Supplementary Tables 10 and 11). The *Rhinacovirus* lineage (L1) that includes HKU2- and SARS-  
160 CoV likely originated in the SO region while all other  $\alpha$ -CoV lineages historically arose in SW China and  
161 spread to other regions before several dispersal events from SO and NO in all directions (Fig. 6A and  
162 Supplementary Fig. 12). A roughly similar pattern of  $\alpha$ -CoV dispersal was obtained using the random  
163 subset (Supplementary Tables 10 and 12).

164 The oldest inferred dispersal movements for  $\beta$ -CoVs occurred among the SO and SW regions (Fig. 6B).  
165 The SO region was the likely origin of *Merbecovirus* (Lineage C, including HKU4- and HKU5-CoV) and  
166 *Sarbecovirus* subgenera (Lineage B, including HKU3- and SARSr-CoVs) while the *Nobecovirus* (lineage D,  
167 including HKU9-CoV) and *Hibecovirus* (lineage E) subgenera originated in SW China (Supplementary Fig.  
168 12). Then several dispersal movements likely originated from SO and CE (Fig. 6B). More recent  
169 southward dispersal from NO was observed. Similar spatiotemporal dispersal patterns were observed  
170 using the random subset of  $\beta$ -CoVs (Supplementary Tables 11 and 13).

171 The estimated rate of migration events per unit of time along these significant dispersal routes was  
172 more than two times higher for  $\alpha$ - (0.026 host switches/unit time) than  $\beta$ -CoVs (0.011 host switches/unit  
173 time) and SO was the region involved in the greatest total number of migration events for both  $\alpha$ - and  $\beta$ -  
174 CoVs. SO had the highest number of outbound and inbound migration events for  $\alpha$ -CoVs (Fig. 6C and  
175 Supplementary Table 12). For  $\beta$ -CoVs, the highest number of outbound migration events was estimated  
176 to be from NO and SO while SO and SW had the highest numbers of inbound migration events (Fig. 6D  
177 and Supplementary Table 13).

### 178 **Phylogenetic diversity**

179 In order to identify the hotspots of CoV phylogenetic diversity in China and evaluate phylogenetic  
180 clustering of CoVs, we calculated the Mean Phylogenetic Distance (MPD) and the Mean Nearest Taxon  
181 Distance (MNTD) statistics<sup>50</sup> and their standardized effect size (SES).

182 We found significant and negative SES MPD values, indicating significant phylogenetic clustering, within  
183 all bat families and genera for both  $\alpha$ - and  $\beta$ -CoVs, except within the *Aselliscus* and *Tylosycteris* for  $\alpha$ -  
184 CoVs (Fig. 7A and 7B). Negative and mostly significant SES MNTD values, reflecting phylogenetic  
185 structure closer to the tips, were also observed within most bat families and genera for  $\alpha$ - and  $\beta$ -CoVs  
186 but we found non-significant positive SES MNTD value for vespertilionid bats, and particularly for those  
187 in the *Pipistrellus* genus, for  $\beta$ -CoVs (Fig. 7A and 7B). In general, we observed lower phylogenetic  
188 diversity for  $\beta$ - than  $\alpha$ -CoVs within all bat families and most genera when looking at SES MPD, but the  
189 difference in the level of diversity between  $\alpha$ - and  $\beta$ -CoVs is less important when looking at SES MNTD  
190 (Fig. 7). These results suggest stronger basal clustering (reflected by larger SES MPD values) for  $\beta$ -CoVs  
191 than  $\alpha$ -CoVs, indicating stronger host structuring effect and phylogenetic conservatism for  $\beta$ -CoVs. Very  
192 similar results were obtained with the random subsets for both  $\alpha$ - and  $\beta$ -CoVs (Supplementary Tables  
193 14-21).

194 We found negative and mostly significant values of MPD and MNTD (Fig. 7C and Supplementary Tables  
195 22-25) indicating significant phylogenetic clustering of CoV lineages in bat communities within the same  
196 zoogeographic region. However, SES MPD values for  $\alpha$ -CoVs in SW were positive (significant for the  
197 random subset) indicating a greater evolutionary diversity of CoVs in that region than others (Fig. 7 and  
198 Supplementary Tables 22-25). We used a linear regression analysis to assess the relationship between  
199 CoV phylogenetic diversity and bat species richness in China and determine if bat richness is a significant  
200 predictor of bat-CoV diversity and evolution.  $\alpha$ -CoV phylogenetic diversity (MPD) was not significantly  
201 correlated to total bat species richness or sampled bat species richness in zoogeographic regions or  
202 provinces (Supplementary Table 26). Non-significant correlations between bat species richness and  $\beta$ -  
203 CoV phylogenetic diversity were also observed at the zoogeographic region level (Supplementary Table  
204 27). However, a significant correlation was observed between sampled bat species richness and  $\beta$ -CoV  
205 phylogenetic diversity at the province level (Supplementary Table 27). Similar results were obtained  
206 when using the random subsets (Supplementary Tables 26 and 27). These findings suggest that bat host  
207 diversity is not the main driver of CoV diversity in China and that other ecological or biogeographic  
208 factors may influence this diversity. We observed higher CoV diversity than expected in several southern  
209 or central provinces (Hainan, Guangxi, Hunan) given their underlying total or sampled bat diversity  
210 (Supplementary Fig. 13 and 14).

211 We also assessed patterns of CoV phylogenetic turnover/differentiation among Chinese zoogeographic  
212 regions and bat host families by measuring the inter-region and inter-host values of MPD (equivalent to  
213 a measure of phylogenetic  $\beta$ -diversity) and their SES. We found positive inter-family SES MPD values,  
214 except between Pteropodidae and Hipposideridae for  $\alpha$ -CoVs and between Rhinolophidae and  
215 Hipposideridae for  $\beta$ -CoVs (Fig. 8A and 8B and Supplementary Tables 28 and 29), suggesting higher  
216 phylogenetic differentiation of CoVs among most bat families than among random communities. Our  
217 phylo-ordination based on inter-family MPD values indicated that  $\alpha$ -CoVs from vespertilionids and

218 miniopterids, and from hipposiderids and pteropodids; as well as  $\beta$ -CoVs from rhinolophids and  
219 hipposiderids are phylogenetically closely related (Fig. 8A and 8B). We also observed strong  
220 phylogenetic turnover between  $\alpha$ -CoV strains from rhinolophids and from miniopterids and all other bat  
221 families, and between  $\beta$ -CoV strains from vespertilionids and all other bat families (Supplementary  
222 Tables 28 and 29). Phylo-ordination among bat genera based on inter-genus MPD confirmed these  
223 results and indicated that CoV strains from genera belonging to the same bat family were mostly more  
224 closely related to each other than to genera from other families (Fig. 8C and 8D and Supplementary  
225 Tables 30 and 31).

226 We observed high and positive inter-region SES MPD values between SW/HI and all other regions,  
227 suggesting that these two regions host higher endemic diversity (Fig. 9 and Supplementary Tables 32  
228 and 31). Negative inter-region SES MPD values suggested that the phylogenetic turnover among other  
229 regions was less important than expected among random communities. Our phylo-ordination among  
230 zoogeographic regions also reflected the high phylogenetic turnover and deep evolutionary  
231 distinctiveness of both  $\alpha$ - and  $\beta$ -CoVs from SW and HI regions (Fig. 9 and Supplementary Tables 32 and  
232 33). Similar results were obtained using the random subset (Supplementary Tables 32 and 33).

### 233 **Mantel tests**

234 Mantel tests revealed a positive and significant correlation between CoV genetic differentiation ( $F_{ST}$ ) and  
235 geographic distance matrices, both with and without provinces including fewer than four viral  
236 sequences, for  $\alpha$ - ( $r = 0.25$ ,  $p = 0.0097$ ;  $r = 0.32$ ,  $p = 0.0196$ ; respectively) and  $\beta$ -CoVs ( $r = 0.22$ ,  $p =$   
237  $0.0095$ ;  $r = 0.23$ ,  $p = 0.0336$ ; respectively). We also detected a positive and highly significant correlation  
238 between CoV genetic differentiation ( $F_{ST}$ ) and their host phylogenetic distance matrices, both with and  
239 without genera including fewer than four viral sequences, for  $\beta$ -CoVs ( $r = 0.41$ ,  $p = 0$ ;  $r = 0.39$ ,  $p =$   
240  $0.0012$ ; respectively) but not for  $\alpha$ -CoVs ( $r = -0.13$ ,  $p = 0.8413$ ;  $r = 0.02$ ,  $p = 0.5019$ ; respectively).

## 241 Discussion

242 Our phylogenetic analysis shows a high diversity of CoVs from bats sampled in China, with most bat  
243 genera included in this study (10/16) infected by both  $\alpha$ - and  $\beta$ -CoVs. In our phylogenetic analysis that  
244 includes all known bat-CoVs from China, we find that SARS-CoV-2 is likely derived from a clade of viruses  
245 originating in horseshoe bats (*Rhinolophus* spp.). The geographic location of this origin appears to be  
246 Yunnan province. However, it is important to note that: 1) our study collected and analyzed samples  
247 solely from China; 2) many sampling sites were close to the borders of Myanmar and Lao PDR; and 3)  
248 most of the bats sampled in Yunnan also occur in these countries, including *R. affinis* and *R. malayanus*,  
249 the species harboring the CoVs with highest *RdRp* sequence identity to SARS-CoV-2<sup>20,21</sup>. For these  
250 reasons, we cannot rule out an origin for the clade of viruses that are progenitors of SARS-CoV-2 that is  
251 outside China, and within Myanmar, Lao PDR, Vietnam or another Southeast Asian country. Additionally,  
252 our analysis shows that the virus RmYN02 from *R. malayanus*, which is characterized by the insertion of  
253 multiple amino acids at the junction site of the S1 and S2 subunits of the Spike (S) protein, belongs to  
254 the same clade as both RaTG13 and SARS-CoV-2, providing further support for the natural origin of  
255 SARS-CoV-2 in *Rhinolophus* spp. bats in the region<sup>20,21</sup>. Finally, while our analysis shows that the *RdRp*  
256 sequences of coronaviruses from the Malayan pangolin are closely related to SARS-CoV-2 *RdRp*, analysis  
257 of full genomes of these viruses suggest that these terrestrial mammals are less likely to be the origin of  
258 SARS-CoV-2 than *Rhinolophus* spp. bats<sup>22,23</sup>.

259 This analysis also demonstrates that a significant amount of cross-species transmission has occurred  
260 among bat hosts over evolutionary time. Our Bayesian phylogeographic inference and analysis of host  
261 switching showed varying levels of viral connectivity among bat hosts and allowed us to identify  
262 significant host transitions that appear to have occurred during bat-CoV evolution in China.

263 We found that bats in the family Rhinolophidae (horseshoe bats) played a key role in the evolution and  
264 cross-species transmission history of  $\alpha$ -CoVs. The family Rhinolophidae and the genus *Rhinolophus* were  
265 involved in more inter-family and inter-genus highly significant host switching of  $\alpha$ -CoVs than any other  
266 family or genus. They were the greatest receivers of  $\alpha$ -CoV host switching events and second greatest  
267 donors after Miniopteridae/*Miniopterus*. The Rhinolophidae, together with the Hipposideridae, also  
268 played an important role in the evolution of  $\beta$ -CoVs, being at the origin of most inter-family host  
269 switching events. Chinese horseshoe bats are characterized by a distinct and evolutionary divergent  $\alpha$ -  
270 CoV diversity, while their  $\beta$ -CoV diversity is similar to that found in the Hipposideridae. The  
271 Rhinolophidae comprises a single genus, *Rhinolophus*, and is the most speciose bat family after the  
272 Vespertilionidae in China<sup>51</sup>, with 20 known species, just under a third of global *Rhinolophus* diversity,  
273 mostly in Southern China<sup>35</sup>. This family likely originated in Asia<sup>52,53</sup>, but some studies suggest an African  
274 origin<sup>54,55</sup>. Rhinolophid fossils from the middle Eocene (38 - 47.8 Mya) have been found in China,  
275 suggesting a westward dispersal of the group from eastern Asia to Europe<sup>56</sup>. The ancient likely origin of  
276 the Rhinolophidae in Asia and China in particular may explain the central role they played in the  
277 evolution and diversification of bat-CoVs in this region, including SARSr-CoVs, MERS-cluster CoVs, and  
278 SADSr-CoVs, which contain important human and livestock pathogens. Horseshoe bats are known to  
279 share roosts with genera from all other bat families in this study<sup>57</sup>, which may also favor CoV cross-  
280 species transmission from and to rhinolophids<sup>34</sup>. A global meta-analysis showing higher rates of viral  
281 sharing among co-roosting cave bats supports this finding<sup>58</sup>.

282 Vespertilionid and miniopterid bats (largely within the *Myotis* and *Miniopterus* genera) also appear to  
283 have been involved in several significant host switches during  $\alpha$ -CoV evolution. However, no significant  
284 transition from vespertilionid bats was identified for  $\beta$ -CoVs and these bats exhibit a divergent  $\beta$ -CoV  
285 diversity compared to other bat families. Vespertilionid and miniopterid bats are characterized by strong  
286 basal phylogenetic clustering but high recent CoV diversification rates, indicating a more rapid

287 evolutionary radiation of CoVs in these bat hosts. At the genus level, similar findings were observed for  
288 the genera *Myotis*, *Pipistrellus* and *Miniopterus*.

289 A significant correlation between geographic distance and genetic differentiation of both  $\alpha$ - and  $\beta$ -CoVs  
290 has been detected, even if only a relatively small proportion of the variance is explained by geographic  
291 distance. We also revealed a significant effect of host phylogeny on  $\beta$ -CoV evolution while it had a  
292 minimal effect on  $\alpha$ -CoV diversity. Contrary to the  $\alpha$ -CoV phylogeny, the basal phylogenetic structure of  
293  $\beta$ -CoVs mirrored the phylogeny of their bat hosts, with a clear distinction between the Yangochiroptera,  
294 encompassing the Vespertilionidae and Miniopteridae, and the Yinpterochiroptera, which includes the  
295 megabat family Pteropodidae and the microbat families Rhinolophidae and Hipposideridae, as  
296 evidenced in recent bat phylogenies<sup>52,59</sup>. These findings suggest a profound co-macroevolutionary  
297 process between  $\beta$ -CoVs and their bat hosts, even if host switches also occurred throughout their  
298 evolution as our study showed. The phylogenetic structure of  $\alpha$ -CoVs, with numerous and closely related  
299 lineages identified in the Vespertilionidae and Miniopteridae, contrasts with the  $\beta$ -CoV  
300 macroevolutionary pattern and suggests  $\alpha$ -CoVs have undergone an adaptive radiation in these two  
301 Yangochiroptera families. Our BSSVS procedure and Markov jump estimates revealed higher  
302 connectivity, both qualitatively and quantitatively, among bat families and genera in the  $\alpha$ -CoV cross-  
303 species transmission history. Larger numbers of highly significant host transitions and higher rates of  
304 switching events along these pathways were inferred for  $\alpha$ - than  $\beta$ -CoVs, especially at the host family  
305 level. These findings suggest that  $\alpha$ -CoVs are able to switch hosts more frequently and between more  
306 distantly related taxa, and that phylogenetic distance among hosts represents a higher constraint on  
307 host switches for  $\beta$ - than  $\alpha$ -CoVs. This is supported by more frequent dispersal events in the evolution of  
308  $\alpha$ - than  $\beta$ -CoVs in China.

309 Variation in the extent of host jumps between  $\alpha$  and  $\beta$ -CoVs within the same hosts in the same  
310 environment may be due to virus-specific factors such as differences in receptor usage between  $\alpha$ - and

311  $\beta$ -CoVs<sup>60-62</sup>. Coronaviruses use a large diversity of receptors, and their entry into host cells is mediated  
312 by the spike protein with an ectodomain consisting of a receptor-binding subunit S1 and a membrane-  
313 fusion subunit S2<sup>63</sup>. However, despite differences in the core structure of their S1 receptor binding  
314 domains (RBD), several  $\alpha$ - and  $\beta$ -CoV species are able to recognize and bind to the same host  
315 receptors<sup>64</sup>. Other factors such as mutation rate, recombination potential, or replication rate might also  
316 be involved in differences in host switching potential between  $\alpha$ - and  $\beta$ -CoVs. A better understanding of  
317 receptor usage and other biological characteristics of these bat-CoVs may help predict their cross-  
318 species transmission and zoonotic potential.

319 We also found that some bat genera were infected by a single CoV genus: *Miniopterus* (Miniopteridae)  
320 and *Murina* (Vespertilionidae) carried only  $\alpha$ -CoVs, while *Cynopterus*, *Eonycteris*, *Megaerops*  
321 (Pteropodidae) and *Pipistrellus* (Vespertilionidae) hosted only  $\beta$ -CoVs. This was found despite using the  
322 same conserved pan-CoV PCR assays for all specimens screened and it can't be explained by differences  
323 in sampling effort for these genera (Supplementary Table 1): for example, >250  $\alpha$ -CoV sequences but no  
324  $\beta$ -CoV were discovered in *Miniopterus* bats in China during our recent fieldwork. These migratory bats,  
325 which seem to have played a key role in the evolution of  $\alpha$ -CoVs, share roosts with several other bat  
326 genera hosting  $\beta$ -CoVs in China<sup>57</sup>, suggesting high likelihood of being exposed to  $\beta$ -CoVs. Biological or  
327 ecological properties of miniopterid bats may explain this observation and clearly warrant further  
328 investigation.

329 Our Bayesian ancestral reconstructions revealed the importance of South western and Southern China  
330 as centers of diversification for both  $\alpha$ - and  $\beta$ -CoVs. These two regions are hotspots of CoV phylogenetic  
331 diversity, harboring evolutionarily old and phylogenetically diverse lineages of  $\alpha$ - and  $\beta$ -CoVs. South  
332 western China acted as a refugium during Quaternary glaciation for numerous plant and animal species  
333 including several bat species, such as *Rhinolophus affinis*<sup>65</sup>, *Rhinolophus sinicus*<sup>66</sup>, *Myotis davidii*<sup>67</sup>, and  
334 *Cynopterus sphinx*<sup>68</sup>. The stable and long-term persistence of bats and other mammals throughout the

335 Quaternary may explain the deep macroevolutionary diversity of bat-CoVs in these regions<sup>69</sup>. Several  
336 highly significant and ancient CoV dispersal routes from these two regions have been identified in this  
337 study. Other viruses, such as the Avian Influenza A viruses H5N6, H7N9 and H5N1, also likely originated  
338 in South western and Southern Chinese regions<sup>70,71</sup>.

339 Our findings suggest that bat host diversity is not the main driver of CoV diversity in China and that  
340 other ecological or biogeographic factors may influence this diversity. Overall, there were no significant  
341 correlations between CoV phylogenetic diversity and bat species diversity (total or sampled) for each  
342 province or biogeographic region, apart from a weak correlation between  $\beta$ -CoV phylogenetic diversity  
343 and the number of bat species sampled at the province level. Yet, we observed higher than expected  
344 phylogenetic diversity in several southern provinces (Hainan, Guangxi, Hunan). These results and main  
345 conclusions are consistent and robust even when we account for geographic biases in sampling effort by  
346 analyzing random subsets of the data.

347 Despite being the most exhaustive study of bat-CoVs in China, this study had several limitations that  
348 must be taken into consideration when interpreting our results. First, only partial *RdRp* sequences were  
349 generated in this study and used in our phylogenetic analysis as the non-invasive samples (rectal  
350 swabs/feces) collected in this study prevented us from generating longer sequences in many cases. The  
351 *RdRp* gene is a suitable marker for this kind of study as it reflects vertical ancestry and is less prone to  
352 recombination than other regions of the CoV genome such as the spike protein gene<sup>16,72</sup>. While using  
353 long sequences is always preferable, our phylogenetic trees are well supported and their topology  
354 consistent with trees obtained using longer sequences or whole genomes<sup>30,73</sup>. Second, most sequences  
355 in this study were obtained by consensus PCR using primers targeting highly conserved regions. Even if  
356 this broadly reactive PCR assay designed to detect widely variant CoVs has proven its ability to detect a  
357 large diversity of CoVs in a wide diversity of bats and mammals<sup>30,74-77</sup>, we may not rule out that some

358 bat-CoV variants remained undetected. Using deep sequencing techniques would allow to detect this  
359 unknown and highly divergent diversity.

360 In this study, we identified the host taxa and geographic regions that together define hotspots of CoV  
361 phylogenetic diversity and centers of diversification in China. These findings may provide a strategy for  
362 targeted discovery of bat-borne CoVs of zoonotic or livestock infection potential, and for early detection  
363 of bat-CoV outbreaks in livestock and people, as proposed elsewhere<sup>78</sup>. Our results suggest that future  
364 sampling and viral discovery should target two hotspots of CoV diversification in Southern and South  
365 western China in particular, as well as neighboring countries where similar bat species live. These  
366 regions are characterized by a subtropical to tropical climate; dense, growing and rapidly urbanizing  
367 populations of people; a high degree of poultry and livestock production; and other factors which may  
368 promote cross-species transmission and disease emergence<sup>78-80</sup>. Additionally, faster rates of evolution in  
369 the tropics have been described for other RNA viruses which could favor cross-species transmission of  
370 RNA viruses in these regions<sup>81</sup>. Both SARS-CoV and SADS-CoV emerged in this region, and several bat  
371 SARSr-CoVs with high zoonotic potential have recently been reported from there, although the dynamics  
372 of their circulation in wild bat populations remain poorly understood<sup>16,61</sup>. Importantly, the closest known  
373 relative of SARS-CoV-2, a SARS-related virus, was found in a *Rhinolophus* sp. bat in this region<sup>20</sup>,  
374 although it is important to note that our survey was limited to China, and that the bat hosts of this virus  
375 also occur in nearby Myanmar and Lao PDR. The significant public health and food security implications  
376 of these outbreaks reinforces the need for enhanced, targeted sampling and discovery of novel CoVs.  
377 Because intensive sampling has not, to our knowledge, been undertaken in countries bordering  
378 southern China, these surveys should be extended to include Myanmar, Lao PDR, and Vietnam, and  
379 perhaps across southeast Asia. Our finding that *Rhinolophus* spp. are most likely to be involved in host-  
380 switching events makes them a key target for future longitudinal surveillance programs, but surveillance

381 targeted the genera *Hipposideros* and *Aselliscus* may also be fruitful as they share numerous  $\beta$ -CoVs  
382 with *Rhinolophus* bats.

383 In the aftermath of the SARS-CoV and MERS-CoV outbreaks,  $\beta$ -CoVs have been the main focus of bat-  
384 CoV studies in China, Africa, and Europe<sup>17,32,36,61,82</sup>. However, we have shown that  $\alpha$ -CoVs have a higher  
385 propensity to switch host within their natural bat reservoirs, and therefore also have a high cross-  
386 species transmission potential and risk of spillover. This is exemplified by the recent emergence of SARS-  
387 CoV in pigs in Guangdong province<sup>17</sup>. Two human  $\alpha$ -CoVs, NL63 and 229E, also likely originated in  
388 bats<sup>27,28</sup>, reminding us that past spillover events from bat species can readily be established in the  
389 human population. Future work discovering and characterizing the biological properties of bat  $\alpha$ -CoVs  
390 may therefore be of potential value for public and livestock health. Our study, and recent analysis of  
391 viral discovery rates<sup>83</sup>, suggest that a substantially wider sampling and discovery net will be required to  
392 capture the complete diversity of coronaviruses in their natural hosts and assess their potential for  
393 cross-species transmission. The bat genera *Rhinolophus*, *Hipposideros*, *Myotis* and *Miniopterus*, all  
394 involved in numerous naturally-occurring host switches throughout  $\alpha$ -CoV evolution, should be a  
395 particular target for  $\alpha$ -CoV discovery in China and across southeast Asia, with *in vitro* and experimental  
396 characterization to better understand their potential to infect people or livestock and cause disease.

## 397 **Methods**

### 398 **Bat sampling**

399 Bat oral and rectal swabs and fecal pellets were collected from 2010 to 2015 in numerous Chinese  
400 provinces (Anhui, Beijing, Guangdong, Guangxi, Guizhou, Hainan, Henan, Hubei, Hunan, Jiangxi, Macau,  
401 Shanxi, Sichuan, Yunnan, and Zhejiang). Fecal pellets were collected from tarps placed below bat  
402 colonies. Bats were captured using mist nets at their roost site or feeding areas. Each captured bat was  
403 stored into a cotton bag, all sampling was non-lethal and bats were released at the site of capture

404 immediately after sample collection. A wing punch was also collected for barcoding purpose. Bat-  
405 handling methods were approved by Tufts University IACUC committee (proposal #G2017-32) and  
406 Wuhan Institute of Virology Chinese Academy of Sciences IACUC committee (proposal WIVA05201705).  
407 Samples were stored in viral transport medium at -80°C directly after collection.

#### 408 **RNA extraction and PCR screening**

409 RNA was extracted from 200 µl swab rectal samples or fecal pellets with the High Pure Viral RNA Kit  
410 (Roche) following the manufacturer's instructions. RNA was eluted in 50 µl elution buffer and stored at -  
411 80°C. A one-step hemi-nested RT-PCR (Invitrogen) was used to detect coronavirus RNA using a set of  
412 primers targeting a 440-nt fragment of the *RdRp* gene and optimized for bat-CoV detection (CoV-FWD3:  
413 GGTTGGGAYTAYCCHAARTGTGA; CoV-RVS3: CCATCATCASWYRAATCATCATA; CoV-FWD4/Bat:  
414 GAYTAYCCHAARTGTGAYAGAGC)<sup>84</sup>. For the first round PCR, the amplification was performed as follows:  
415 50°C for 30 min, 94°C for 2 min, followed by 40 cycles consisting of 94°C for 20 sec, 50°C for 30 sec, 68°C  
416 for 30 sec, and a final extension step at 68°C for 5 min. For the second round PCR, the amplification was  
417 performed as follows: 94°C for 2 min followed by 40 cycles consisting of 94°C for 20 sec, 59°C for 30 sec,  
418 72°C for 30 sec, and a final extension step at 72°C for 7 min. PCR products were gel purified and  
419 sequenced with an ABI Prism 3730 DNA analyzer (Applied Biosystems, USA). PCR products with low  
420 concentration or bad sequencing quality were cloned into pGEM-T Easy Vector (Promega) for  
421 sequencing. Positive results detected in bat genera that were not known to harbor a specific CoV lineage  
422 previously were repeated a second time (PCR + sequencing) as a confirmation. Species identifications  
423 from the field were also confirmed and re-confirmed by cytochrome (cytb) DNA barcoding using DNA  
424 extracted from the feces or swabs<sup>85</sup>. Only viral detection and barcoding results confirmed at least twice  
425 were included in this study.

#### 426 **Sequence data**

427 We also added bat-CoV *RdRp* sequences from China available in GenBank to our dataset. All sequences  
428 for which sampling year and host or sampling location information was available either in GenBank  
429 metadata or in the original publication were included (as of March 15, 2018). Our final datasets include  
430 630 sequences generated for this study and 616 sequences from GenBank or GISAID (list of GenBank  
431 and GISAID accession numbers available in Supplementary Note 1, and Supplementary Tables 34 and  
432 35). Nucleotide sequences were aligned using MUSCLE and trimmed to 360 base pair length to reduce  
433 the proportion of missing data in the alignments. All phylogenetic analyses were performed on both the  
434 complete data and random subset, and for  $\alpha$ - and  $\beta$ -CoVs separately.

#### 435 **Defining zoogeographic regions in China**

436 Hierarchical clustering was used to define zoogeographic regions within China by clustering provinces  
437 with similar mammalian diversity<sup>45</sup>. Hierarchical cluster analysis classifies several objects into small  
438 groups based on similarities between them. To do this, we created a presence/absence matrix of all  
439 extant terrestrial mammals present in China using data from the IUCN spatial database<sup>86</sup> and generated  
440 a cluster dendrogram using the function *hclust* with average method of the R package *stats*. Hong Kong  
441 and Macau were included within the neighboring Guangdong province. We then visually identified  
442 geographically contiguous clusters of provinces for which CoV sequences are available (Fig. 1 and  
443 Supplementary Fig. 1).

444 We identified six zoogeographic regions within China based on the similarity of the mammal community  
445 in these provinces: South western region (SW; Yunnan province), Northern region (NO; Xizang, Gansu,  
446 Jilin, Anhui, Henan, Shandong, Shaanxi, Hebei and Shanxi provinces and Beijing municipality), Central  
447 northern region (CN; Sichuan and Hubei provinces), Central region (CE; Guangxi, Guizhou, Hunan, Jiangxi  
448 and Zhejiang provinces), Southern region (SO; Guangdong and Fujian provinces, Hong Kong, Macau and  
449 Taiwan), and Hainan island (HI). Hunan and Jiangxi, clustering with the SO provinces in our dendrogram,

450 were included within the central region to create a geographically contiguous Central cluster  
451 (Supplementary Fig. 1). These six zoogeographic regions are very similar to the biogeographic regions  
452 traditionally recognized in China<sup>87</sup>. The three  $\beta$ -CoV sequences from HI were included in the SO region to  
453 avoid creating a cluster with a very small number of sequences.

#### 454 **Model selection and phylogenetic analysis**

455 Bayesian phylogenetic analysis were performed in BEAST 1.8.4<sup>46</sup>. Sampling years were used as tip dates.  
456 Preliminary analysis were run to select the best fitting combination of substitution models (HKY/GTR),  
457 codon partition scheme, molecular clock (strict/lognormal uncorrelated relaxed clock) and coalescent  
458 models (constant population size/exponential growth/GMRF Bayesian Skyride). Model combinations  
459 were compared and the best fitting model was selected using a modified Akaike information criterion  
460 (AICM) implemented in Tracer 1.6<sup>88</sup>. We also used TEMPEST<sup>89</sup> to assess the temporal structure within  
461 our  $\alpha$ - and  $\beta$ -CoV datasets. TEMPEST showed that both datasets did not contain sufficient temporal  
462 information to accurately estimate substitution rates or time to the most recent common ancestor  
463 (TMRCA). Therefore we used a fixed substitution rate of 1.0 for all our BEAST analysis.

464 All subsequent BEAST analysis were performed under the best fitting model including a HKY substitution  
465 model with two codons partitions ((1+2), 3), a strict molecular clock and a constant population size  
466 coalescent model. Each analysis was run for  $2.5 \times 10^8$  generations, with sampling every  $2 \times 10^4$  steps. All  
467 BEAST computations were performed on the CIPRES Science Getaway Portal<sup>90</sup>. Convergence of the chain  
468 was assessed in Tracer so that the effective sample size (ESS) of all parameters was  $> 200$  after removing  
469 at least 10% of the chain as burn-in.

#### 470 **Ancestral state reconstruction and transition rates**

471 A Bayesian discrete phylogeographic approach implemented in BEAST 1.8.4 was used to reconstruct the  
472 ancestral state of each node in the phylogenetic tree for three discrete traits: host family, host genus

473 and zoogeographic region. An asymmetric trait substitution model was applied. These analyses were  
474 performed for each trait on the complete dataset and random subsets. Maximum clade credibility (MCC)  
475 tree annotated with discrete traits were generated in TreeAnnotator and visualized using the software  
476 Spred3<sup>91</sup>.

477 For each analysis, a Bayesian stochastic search variable selection (BSSVS) was applied to estimate the  
478 significance of pairwise switches between trait states using Bayesian Factor (BF) as a measure of  
479 statistical significance<sup>47</sup>. BF were computed in Spred3. BF support was interpreted according to Jeffreys  
480 1961<sup>92</sup> (BF > 3: substantial support, BF > 10: strong support, BF > 30: very strong support, BF > 100:  
481 decisive support) and only strongly supported transitions were presented in most figures, following a  
482 strategy used in other studies<sup>93,94</sup>. We also estimated the count of state switching events (Markov  
483 jumps)<sup>48,49</sup> along the branches of the phylogenetic tree globally (for the three discrete traits) and for  
484 each strongly supported (BF > 10) transition between character states (for bat families and ecoregions  
485 only). Convergence of the MCMC runs was confirmed using Tracer. The rate of state switching events  
486 per unit of time was estimated for each CoV genus by dividing the total estimated number of state  
487 switching events by the total branch length of the MCC tree.

488 To assess the phylogenetic relationships among SARS-CoV-2 and other CoVs from the *Sarbecovirus*  
489 subgenus, we also reconstructed a MCC tree in BEAST 1.8.4 and median-joining network in Network  
490 10.0<sup>95</sup> including all *Sarbecovirus* sequences, two sequences of SARS-CoV-2 isolated in humans (GenBank  
491 accession numbers: MN908947 and MN975262), one sequence of SARS-CoV (GenBank accession  
492 number: NC\_004718), eight sequences from Malayan pangolins (*Manis javanica*) (GISAID accession  
493 numbers: EPI\_ISL\_410538-410544, EPI\_ISL\_410721) and one from *Rhinolophus malayanus* (GISAID  
494 accession number: EPI\_ISL\_412977).

#### 495 **Phylogenetic diversity**

496 The Mean Phylogenetic Distance (MPD) and the Mean Nearest Taxon Distance (MNTD) statistics<sup>50</sup> and  
497 their standardized effect size (SES) were calculated for each zoogeographic region, bat family and genus  
498 using the R package *picante*<sup>96</sup>. MPD measures the mean phylogenetic distance among all pairs of CoVs  
499 within a host or a region. It reflects phylogenetic structuring across the whole phylogenetic tree and  
500 assesses the overall divergence of CoV lineages in a community. MNTD is the mean distance between  
501 each CoV and its nearest phylogenetic neighbor in a host or region, and therefore it reflects the  
502 phylogenetic structuring closer to the tips and shows how locally clustered taxa are. SES MPD and SES  
503 MNTD values correspond to the difference between the phylogenetic distances in the observed  
504 communities versus null communities. Low and negative SES values denote phylogenetic clustering, high  
505 and positive values indicate phylogenetic over-dispersion while values close to 0 show random  
506 dispersion. The SES values were calculated by building null communities by randomly reshuffling tip  
507 labels 1000 times along the entire phylogeny. Phylogenetic diversity computations were performed on  
508 both the complete dataset and random subset for each trait. A linear regression analysis was performed  
509 in R to assess the correlation between CoV phylogenetic diversity (MPD) and bat species richness in  
510 China. Total species richness per province or region was estimated using data from the IUCN spatial  
511 database while sampled species richness corresponds to the number of bat species sampled and tested  
512 for CoV per province or region in our datasets.

513 The inter-region and inter-host values of MPD (equivalent to phylogenetic  $\beta$  diversity), corresponding to  
514 the mean phylogenetic distance among all pairs of CoVs from two distinct hosts or regions, and their SES  
515 were estimated using the function *comdist* of the R package *phylocomr*<sup>97</sup>. The matrices of inter-region  
516 and inter-host MPD were used to cluster zoogeographic regions and bat hosts in a dendrogram  
517 according to their evolutionary similarity (phylo-ordination) using the function *hclust* with complete  
518 linkage method of the R package *stats* (R core team). These computations were performed on both the  
519 complete dataset and random subset.

## 520 **Mantel tests and isolation by distance**

521 Mantel tests performed in ARLEQUIN 3.5<sup>98</sup> were used to compare the matrix of viral genetic  
522 differentiation ( $F_{ST}$ ) to matrices of host phylogenetic distance and geographic distance in order to  
523 evaluate the role of geographic isolation and host phylogeny in shaping CoV population structure. The  
524 correlation between these matrices was assessed using 10,000 permutations. To gain more resolution  
525 into the process of evolutionary diversification, these analyses were also performed at the host genus  
526 and province levels. To calculate phylogenetic distances among bat genera, we reconstructed a  
527 phylogenetic tree including a single sequence for all bat species included in our dataset. Pairwise  
528 patristic distances among tips were computed using the function *distTips* in the R package *adephylo*<sup>99</sup>.  
529 We then averaged all distances across genera to create a matrix of pairwise distances among bat  
530 genera. Pairwise Euclidian distances were measured between province centroids and log transformed.  
531 Mantel tests were performed with and without genera and provinces including less than four viral  
532 sequences to assess the impact of low sample size on our results.

## 533 **Data availability**

534 GenBank accession numbers of sequences generated in this study and previously published sequences  
535 included in our analysis are available in the Supplementary Note 1 and Supplementary Tables 34 and 35.

## 536 **References**

- 537 1. Forni, D., Cagliani, R., Clerici, M. & Sironi, M. Molecular Evolution of Human Coronavirus  
538 Genomes. *Trends in Microbiology* **25**, 35-48 (2017).
- 539 2. Tao, Y. *et al.* Surveillance of Bat Coronaviruses in Kenya Identifies Relatives of Human  
540 Coronaviruses NL63 and 229E and Their Recombination History. *Journal of Virology* **91**(2017).

- 541 3. Graham, R.L. & Baric, R.S. Recombination, Reservoirs, and the Modular Spike: Mechanisms of  
542 Coronavirus Cross-Species Transmission. **84**, 3134-3146 (2010).
- 543 4. Vijgen, L. *et al.* Evolutionary history of the closely related group 2 coronaviruses: porcine  
544 hemagglutinating encephalomyelitis virus, bovine coronavirus, and human coronavirus OC43.  
545 *Journal of virology* **80**, 7270-7274 (2006).
- 546 5. Zhang, X. *et al.* Quasispecies of bovine enteric and respiratory coronaviruses based on complete  
547 genome sequences and genetic changes after tissue culture adaptation. *Virology* **363**, 1-10  
548 (2007).
- 549 6. Parrish, C.R. *et al.* Cross-Species Virus Transmission and the Emergence of New Epidemic  
550 Diseases. *Microbiology and Molecular Biology Reviews* **72**, 457-470 (2008).
- 551 7. Li, D.L. *et al.* Molecular evolution of porcine epidemic diarrhea virus and porcine  
552 deltacoronavirus strains in Central China. *Research in Veterinary Science* **120**, 63-69 (2018).
- 553 8. Cui, J., Li, F. & Shi, Z.-L. Origin and evolution of pathogenic coronaviruses. *Nature Reviews*  
554 *Microbiology* **17**, 181-192 (2019).
- 555 9. Lau, S.K.P. & Chan, J.F.W. Coronaviruses: emerging and re-emerging pathogens in humans and  
556 animals. *Virology Journal* **12**, 209 (2015).
- 557 10. Drosten, C. *et al.* Identification of a novel coronavirus in patients with severe acute respiratory  
558 syndrome. *N Engl J Med* **348**, 1967-76 (2003).
- 559 11. Heymann, D.L. The international response to the outbreak of SARS in 2003. *Philosophical*  
560 *Transactions of the Royal Society of London Series B-Biological Sciences* **359**, 1127-1129 (2004).
- 561 12. World Health Organization. Summary of probable SARS cases with onset of illness from 1  
562 November 2002 to 31 July 2003. Vol. 2019 (World Health Organization, 2004).
- 563 13. Ge, X.-Y. *et al.* Isolation and characterization of a bat SARS-like coronavirus that uses the ACE2  
564 receptor. *Nature* **503**, 535-538 (2013).

- 565 14. Li, W. *et al.* Bats are natural reservoirs of SARS-like coronaviruses. *Science* **310**, 676-9 (2005).
- 566 15. Lau, S.K.P. *et al.* Severe acute respiratory syndrome coronavirus-like virus in Chinese horseshoe  
567 bats. *Proceedings of the National Academy of Sciences of the United States of America* **102**,  
568 14040-14045 (2005).
- 569 16. Hu, B. *et al.* Discovery of a rich gene pool of bat SARS-related coronaviruses provides new  
570 insights into the origin of SARS coronavirus. *PLoS Pathogens* **13**, e1006698 (2017).
- 571 17. Zhou, P. *et al.* Fatal swine acute diarrhoea syndrome caused by an HKU2-related coronavirus of  
572 bat origin. *Nature* **556**, 255-258 (2018).
- 573 18. Gong, L. *et al.* A New Bat-HKU2-like Coronavirus in Swine, China, 2017. *Emerging infectious*  
574 *diseases* **23**, 1607-1609 (2017).
- 575 19. Pan, Y. *et al.* Discovery of a novel swine enteric alphacoronavirus (SeACoV) in southern China.  
576 *Veterinary Microbiology* **211**, 15-21 (2017).
- 577 20. Zhou, P., Yang, X.-L., Wang, X.-G., Hu, B., Zhang, L., Zhang, W., Si, H.-R., Zhu, Y., Li, B., Huang, C.-  
578 L., *et al.* A pneumonia outbreak associated with a new coronavirus of probable bat origin.  
579 *Nature*, **579**, 270-273 (2020).
- 580 21. Zhou, H., Chen, X., Hu, T., Li, J., Song, H., Liu, Y., Wang, P., Liu, D., Yang, J., Holmes, E.C., *et al.* A  
581 novel bat coronavirus closely related to SARS-CoV-2 contains natural insertions at the S1/S2  
582 cleavage site of the spike protein. *Current Biology* (2020).
- 583 22. Lam, T.T.-Y., Shum, M.H.-H., Zhu, H.-C., Tong, Y.-G., Ni, X.-B., Liao, Y.-S., Wei, W., Cheung, W.Y.-  
584 M., Li, W.-J., Li, L.-F., *et al.* Identifying SARS-CoV-2 related coronaviruses in Malayan pangolins.  
585 *Nature* (2020).
- 586 23. Xiao, K., Zhai, J., Feng, Y., Zhou, N., Zhang, X., Zou, J.-J., Li, N., Guo, Y., Li, X., Shen, X., *et al.*  
587 Isolation of SARS-CoV-2-related coronavirus from Malayan pangolins. *Nature* (2020).

- 588 24. Corman, V.M. *et al.* Rooting the Phylogenetic Tree of Middle East Respiratory Syndrome  
589 Coronavirus by Characterization of a Conspecific Virus from an African Bat. *Journal of Virology*  
590 **88**, 11297-11303 (2014).
- 591 25. Anthony, S.J. *et al.* Further Evidence for Bats as the Evolutionary Source of Middle East  
592 Respiratory Syndrome Coronavirus. *mBio* **8**, e00373-17 (2017).
- 593 26. Lau, S.K.P. *et al.* Receptor Usage of a Novel Bat Lineage C Betacoronavirus Reveals Evolution of  
594 Middle East Respiratory Syndrome-Related Coronavirus Spike Proteins for Human Dipeptidyl  
595 Peptidase 4 Binding. *The Journal of Infectious Diseases*, jiy018-jiy018 (2018).
- 596 27. Corman, V.M. *et al.* Evidence for an Ancestral Association of Human Coronavirus 229E with Bats.  
597 *Journal of Virology* **89**, 11858-11870 (2015).
- 598 28. Huynh, J. *et al.* Evidence Supporting a Zoonotic Origin of Human Coronavirus Strain NL63.  
599 *Journal of Virology* **86**, 12816-12825 (2012).
- 600 29. Lu, R., Zhao, X., Li, J., Niu, P., Yang, B., Wu, H., Wang, W., Song, H., Huang, B., Zhu, N., *et al.*  
601 Genomic characterisation and epidemiology of 2019 novel coronavirus: implications for virus  
602 origins and receptor binding. *The Lancet* **395**, 565-574 (2020).
- 603 30. Wong, A.C.P., Li, X., Lau, S.K.P. & Woo, P.C.Y. Global Epidemiology of Bat Coronaviruses. *Viruses*  
604 **11**, 174 (2019).
- 605 31. Drexler, J.F., Corman, V.M. & Drosten, C. Ecology, evolution and classification of bat  
606 coronaviruses in the aftermath of SARS. *Antiviral Research* **101**, 45-56 (2014).
- 607 32. Anthony, S.J. *et al.* Global patterns in coronavirus diversity. *Virus Evolution* **3**, vex012-vex012  
608 (2017).
- 609 33. Leopardi, S. *et al.* Interplay between co-divergence and cross-species transmission in the  
610 evolutionary history of bat coronaviruses. *Infection, Genetics and Evolution* **58**, 279-289 (2018).

- 611 34. Cui, J. *et al.* Evolutionary relationships between bat coronaviruses and their hosts. *Emerging*  
612 *Infectious Diseases* **13**, 1526-1532 (2007).
- 613 35. Smith, A.T. & Xie, Y. *A Guide to the Mammals of China*, (Princeton University Press, Princeton,  
614 USA, 2008).
- 615 36. Lin, X.-D. *et al.* Extensive diversity of coronaviruses in bats from China. *Virology* **507**, 1-10 (2017).
- 616 37. Ge, X.-Y. *et al.* Coexistence of multiple coronaviruses in several bat colonies in an abandoned  
617 mineshaft. *Virologica Sinica* **31**, 31-40 (2016).
- 618 38. Woo, P.C.Y. *et al.* Molecular diversity of coronaviruses in bats. *Virology* **351**, 180-187 (2006).
- 619 39. Wu, Z. *et al.* Deciphering the bat virome catalog to better understand the ecological diversity of  
620 bat viruses and the bat origin of emerging infectious diseases. *The Isme Journal* **10**, 609-620  
621 (2016).
- 622 40. Tang, X.C. *et al.* Prevalence and Genetic Diversity of Coronaviruses in Bats from China. *Journal of*  
623 *Virology* **80**, 7481-7490 (2006).
- 624 41. Woo, P.C.Y. *et al.* Comparative Analysis of Twelve Genomes of Three Novel Group 2c and Group  
625 2d Coronaviruses Reveals Unique Group and Subgroup Features. *Journal of Virology* **81**, 1574-  
626 1585 (2007).
- 627 42. Ge, X. *et al.* Metagenomic analysis of viruses from bat fecal samples reveals many novel viruses  
628 in insectivorous bats in China. *J Virol* **86**, 4620-4630 (2012).
- 629 43. Xu, L. *et al.* Detection and characterization of diverse alpha- and betacoronaviruses from bats in  
630 China. *Virologica Sinica* **31**, 69-77 (2016).
- 631 44. Luo, Y. *et al.* Longitudinal Surveillance of Betacoronaviruses in Fruit Bats in Yunnan Province,  
632 China During 2009–2016. **33**, 87-95 (2018).
- 633 45. Legendre, P. & Legendre, L.F. *Numerical ecology*, (Elsevier, 2012).

- 634 46. Drummond, A.J., Suchard, M.A., Xie, D. & Rambaut, A. Bayesian phylogenetics with BEAUti and  
635 the BEAST 1.7. *Molecular Biology and Evolution* **29**, 1969-1973 (2012).
- 636 47. Lemey, P., Rambaut, A., Drummond, A.J. & Suchard, M.A. Bayesian Phylogeography Finds Its  
637 Roots. *PLoS Computational Biology* **5**, e1000520 (2009).
- 638 48. Minin, V.N. & Suchard, M.A. Counting labeled transitions in continuous-time Markov models of  
639 evolution. *Journal of Mathematical Biology* **56**, 391-412 (2008).
- 640 49. O'Brien, J.D., Minin, V.N. & Suchard, M.A. Learning to Count: Robust Estimates for Labeled  
641 Distances between Molecular Sequences. *Molecular Biology and Evolution* **26**, 801-814 (2009).
- 642 50. Webb, C.O., Ackerly, D.D., McPeck, M.A. & Donoghue, M.J. Phylogenies and Community  
643 Ecology. **33**, 475-505 (2002).
- 644 51. Simmons, N.B. Order Chiroptera. in *Mammal Species of the World: A Taxonomic and Geographic*  
645 *Reference* (eds. Wilson, D.E. & Reeder, D.M.) 312-529 (Johns Hopkins University Press, 2005).
- 646 52. Teeling, E.C. *et al.* A Molecular Phylogeny for Bats Illuminates Biogeography and the Fossil  
647 Record. **307**, 580-584 (2005).
- 648 53. Stoffberg, S., Jacobs, D.S., Mackie, I.J. & Matthee, C.A. Molecular phylogenetics and historical  
649 biogeography of *Rhinolophus* bats. *Molecular Phylogenetics and Evolution* **54**, 1-9 (2010).
- 650 54. Foley, N.M. *et al.* How and Why Overcome the Impediments to Resolution: Lessons from  
651 rhinolophid and hipposiderid Bats. *Molecular Biology and Evolution* **32**, 313-333 (2014).
- 652 55. Eick, G.N., Jacobs, D.S. & Matthee, C.A. A Nuclear DNA Phylogenetic Perspective on the  
653 Evolution of Echolocation and Historical Biogeography of Extant Bats (Chiroptera). *Molecular*  
654 *Biology and Evolution* **22**, 1869-1886 (2005).
- 655 56. Ravel, A., Marivaux, L., Qi, T., Wang, Y.-Q. & Beard, K.C. New chiropterans from the middle  
656 Eocene of Shanghuang (Jiangsu Province, Coastal China): new insight into the dawn horseshoe  
657 bats (Rhinolophidae) in Asia. **43**, 1-23 (2014).

- 658 57. Luo, J. *et al.* Bat conservation in China: should protection of subterranean habitats be a priority?  
659 *Oryx* **47**, 526-531 (2013).
- 660 58. Willoughby, A.R., Phelps, K.L., Consortium, P. & Olival, K.J. A Comparative Analysis of Viral  
661 Richness and Viral Sharing in Cave-Roosting Bats. *Diversity* **9**, 35 (2017).
- 662 59. Tsagkogeorga, G., Parker, J., Stupka, E., Cotton, James A. & Rossiter, S.J. Phylogenomic Analyses  
663 Elucidate the Evolutionary Relationships of Bats. *Current Biology* **23**, 2262-2267 (2013).
- 664 60. Yang, Y. *et al.* Receptor usage and cell entry of bat coronavirus HKU4 provide insight into bat-to-  
665 human transmission of MERS coronavirus. *Proceedings of the National Academy of Sciences* **111**,  
666 12516-12521 (2014).
- 667 61. Menachery, V.D. *et al.* A SARS-like cluster of circulating bat coronaviruses shows potential for  
668 human emergence. *Nature Medicine* **21**, 1508-1513 (2015).
- 669 62. Li, W. *et al.* Angiotensin-converting enzyme 2 is a functional receptor for the SARS coronavirus.  
670 *Nature* **426**, 450-454 (2003).
- 671 63. Li, F. Receptor Recognition Mechanisms of Coronaviruses: a Decade of Structural Studies. **89**,  
672 1954-1964 (2015).
- 673 64. Li, F. Structure, Function, and Evolution of Coronavirus Spike Proteins. *Annual Review of Virology*  
674 **3**, 237-261 (2016).
- 675 65. Mao, X.G., Zhu, G.J., Zhang, S. & Rossiter, S.J. Pleistocene climatic cycling drives intra-specific  
676 diversification in the intermediate horseshoe bat (*Rhinolophus affinis*) in Southern China.  
677 *Molecular Ecology* **19**, 2754-2769 (2010).
- 678 66. Mao, X. *et al.* Multiple cases of asymmetric introgression among horseshoe bats detected by  
679 phylogenetic conflicts across loci. *Biological Journal of the Linnean Society* **110**, 346-361 (2013).
- 680 67. You, Y. *et al.* Pleistocene glacial cycle effects on the phylogeography of the Chinese endemic bat  
681 species, *Myotis davidii*. *BMC Evolutionary Biology* **10**, 208 (2010).

- 682 68. Chen, J.P. *et al.* Contrasting Genetic Structure in Two Co-Distributed Species of Old World Fruit  
683 Bat. *PLoS ONE* **5** (2010).
- 684 69. Krasnov, B.R., Piloosof, S., Shenbrot, G.I. & Khokhlova, I.S. Spatial variation in the phylogenetic  
685 structure of flea assemblages across geographic ranges of small mammalian hosts in the  
686 Palearctic. *International Journal for Parasitology* **43**, 763-770 (2013).
- 687 70. Bi, Y. *et al.* Novel avian influenza A (H5N6) viruses isolated in migratory waterfowl before the  
688 first human case reported in China, 2014. *Scientific Reports* **6**, 29888 (2016).
- 689 71. Bui, C.M., Adam, D.C., Njoto, E., Scotch, M. & MacIntyre, C.R. Characterising routes of H5N1 and  
690 H7N9 spread in China using Bayesian phylogeographical analysis. *Emerging Microbes &*  
691 *Infections* **7**, 184 (2018).
- 692 72. Gouilh, M.A., Puechmaille, S.J., Gonzalez, J.-P., Teeling, E., Kittayapong, P. & Manuguerra, J.-C.  
693 SARS-Coronavirus ancestor's foot-prints in South-East Asian bat colonies and the refuge theory.  
694 *Infection Genetics and Evolution* **11**, 1690-1702 (2011).
- 695 73. Hu, B., Ge, X., Wang, L.-F. & Shi, Z. Bat origin of human coronaviruses. *Virology Journal* **12**, 1-10  
696 (2015).
- 697 74. Anthony, S.J., Ojeda-Flores, R., Rico-Chávez, O., Navarrete-Macias, I., Zambrana-Torrel, C.M.,  
698 Rostal, M.K., Epstein, J.H., Tipps, T., Liang, E., Sanchez-Leon, M., *et al.* Coronaviruses in bats from  
699 Mexico. *Journal of General Virology* **94**, 1028-1038 (2013).
- 700 75. Corman, V.M., Kallies, R., Philipps, H., Göpner, G., Müller, M.A., Eckerle, I., Brünink, S., Drosten,  
701 C. & Drexler, J.F. Characterization of a novel betacoronavirus related to MERS-CoV in European  
702 hedgehogs. *Journal of Virology* **88**, 717-724 (2014).
- 703 76. Munster, V.J., Adney, D.R., van Doremalen, N., Brown, V.R., Miazgowiec, K.L., Milne-Price, S.,  
704 Bushmaker, T., Rosenke, R., Scott, D., Hawkinson, A., *et al.* Replication and shedding of MERS-  
705 CoV in Jamaican fruit bats (*Artibeus jamaicensis*). *Scientific Reports* **6**, 21878 (2016).

- 706 77. Joyjinda, Y., Rodpan, A., Chartpituck, P., Suthum, K., Yaemsakul, S., Cheun-Arom, T., Bunprakob,  
707 S., Olival, K.J., Stokes, M.M., Hemachudha, T., *et al.* First Complete Genome Sequence of Human  
708 Coronavirus HKU1 from a Nonill Bat Guano Miner in Thailand. *Microbiology Resource*  
709 *Announcements* **8**, e01457-01418 (2019).
- 710 78. Carroll, D., Daszak, P., Wolfe, N.D., Gao, G.F., Morel, C.M., Morzaria, S., Pablos-Méndez, A.,  
711 Tomori, O. & Mazet, J.A.K. The Global Virome Project. *Science* **359**, 872-874 (2018).
- 712 79. Fountain-Jones, N.M. *et al.* Towards an eco-phylogenetic framework for infectious disease  
713 ecology. **93**, 950-970 (2018).
- 714 80. Allen, T. *et al.* Global hotspots and correlates of emerging zoonotic diseases. *Nature*  
715 *Communications* **8**, 1124 (2017).
- 716 81. Streicker, D.G., Lemey, P., Velasco-Villa, A. & Rupprecht, C.E. Rates of Viral Evolution Are Linked  
717 to Host Geography in Bat Rabies. *PLoS Pathog* **8**, e1002720 (2012).
- 718 82. Hu, B. *et al.* Discovery of a rich gene pool of bat SARS-related coronaviruses provides new  
719 insights into the origin of SARS coronavirus. *PLOS Pathogens* **13**(2017).
- 720 83. Carroll, D. *et al.* The global virome project. *Science* **359**, 872-874 (2018).
- 721 84. Watanabe, S. *et al.* Bat Coronaviruses and Experimental Infection of Bats, the Philippines.  
722 *Emerging Infectious Diseases* **16**, 1217-1223 (2010).
- 723 85. Irwin, D.M., Kocher, T.D. & Wilson, A.C. Evolution of the cytochrome b gene of mammals.  
724 *Journal of Molecular Evolution* **32**, 128-144 (1991).
- 725 86. IUCN. The IUCN Red List of Threatened Species. Version 2015.2, <http://www.iucnredlist.org>.  
726 (2018).
- 727 87. Xie, Y., MacKinnon, J., Li, D.J.B. & Conservation. Study on biogeographical divisions of China. **13**,  
728 1391-1417 (2004).

- 729 88. Baele, G., Li, W.L.S., Drummond, A.J., Suchard, M.A. & Lemey, P. Accurate Model Selection of  
730 Relaxed Molecular Clocks in Bayesian Phylogenetics. *Molecular Biology and Evolution* **30**, 239-  
731 243 (2013).
- 732 89. Rambaut, A., Lam, T.T., Max Carvalho, L. & Pybus, O.G. Exploring the temporal structure of  
733 heterochronous sequences using TempEst (formerly Path-O-Gen). *Virus Evolution* **2**, vew007  
734 (2016).
- 735 90. Miller, M.A., Pfeiffer, W. & Schwartz, T. Creating the CIPRES Science Gateway for inference of  
736 large phylogenetic trees. *Proceedings of the Gateway Computing Environments Workshop (GCE)*,  
737 *14 Nov. 2010, New Orleans, LA*, 1-8 (2010).
- 738 91. Bielejec, F. *et al.* Spred3: Interactive Visualization of Spatiotemporal History and Trait  
739 Evolutionary Processes. *Molecular Biology and Evolution* **33**, 2167-2169 (2016).
- 740 92. Jeffreys, H. *Theory of probability*. Oxford: Clarendon. (1961).
- 741 93. Faria, N.R., Suchard, M.A., Rambaut, A., Streicker, D.G. & Lemey, P. Simultaneously  
742 reconstructing viral cross-species transmission history and identifying the underlying  
743 constraints. *Philosophical Transactions of the Royal Society B: Biological Sciences* **368**, 20120196  
744 (2013).
- 745 94. Kamath, P.L., Foster, J.T., Drees, K.P., Luikart, G., Quance, C., Anderson, N.J., Clarke, P.R., Cole,  
746 E.K., Drew, M.L., Edwards, W.H., *et al.* Genomics reveals historic and contemporary transmission  
747 dynamics of a bacterial disease among wildlife and livestock. *Nature Communications* **7**, 11448  
748 (2016).
- 749 95. Bandelt, H.J., Forster, P., & Rohl, A. Median-joining networks for inferring intraspecific  
750 phylogenies. *Molecular Biology and Evolution* **16**, 37-48 (1999).
- 751 96. Kembel, S.W. *et al.* Picante: R tools for integrating phylogenies and ecology. *Bioinformatics* **26**,  
752 1463-1464 (2010).

- 753 97. Ooms, J., Chamberlain, S., Webb, C.O., Ackerly, D.D. & Kembel, S.W. phylocomr: Interface to  
754 'Phylocom'. *R package version 0.1.2* (2018).
- 755 98. Excoffier, L. & Lischer, H.E.L. Arlequin suite ver 3.5: a new series of programs to perform  
756 population genetics analyses under Linux and Windows. *Molecular Ecology Resources* **10**, 564-  
757 567 (2010).
- 758 99. Jombart, T. & Dray, S. adephylo: exploratory analyses for the phylogenetic comparative method.  
759 R package version 1.1-11. (2008).

## 760 **Acknowledgements**

761 This study was funded by the National Institute of Allergy and Infectious Diseases of the National  
762 Institutes of Health (Award Number R01AI110964) and the United States Agency for International  
763 Development (USAID) Emerging Pandemic Threats PREDICT project (cooperative agreement number  
764 GHN-A-OO-09-00010-00), the strategic priority research program of the Chinese Academy of Sciences  
765 (XDB29010101), and National Natural Science Foundation of China (31770175, 31830096). Coronavirus  
766 research in L-FW's group is funded by grants from Singapore National Research Foundation  
767 (NRF2012NRF-CRP001-056 and NRF2016NRF-NSFC002-013).

## 768 **Author contributions**

769 K.J.O., H.E.F, J.H.E., L-F.W., Z.S. and P.D. created the study design, initiated field work and set up sample  
770 collection and testing protocols. B.H., G.Z., L.Z., H.L., A.A.C and Z.L. collected samples or provided  
771 data. B.H., B.L., and W.Z. performed laboratory work. A.L. carried out the analyses and drafted the  
772 manuscript with K.J.O, C.Z.-T. and P.D. All authors reviewed and edited the manuscript

773 **Competing interests:** The authors declare no competing interests.

## 774 **Figure legends**

775 **Fig. 1 Geographic sampling.** Pie chart (A) showing the number of sequences of each CoV genus (alpha-  
776 CoVs and beta-CoVs) available for each zoogeographic region and map of China provinces (B) showing  
777 the number of *RdRp* sequences available for each province, in bold grey for alpha-CoVs and black for  
778 beta-CoVs. Province colors correspond to the zoogeographic region to which they belong: NO, Northern  
779 region; CN, Central northern region; SW, South western region; CE, Central region; SO, Southern region;  
780 HI, Hainan island. The three beta-CoV sequences from HI were included in the SO region. Provinces  
781 colored in grey are those where CoV sequences are not available.

782 **Fig. 2 Phylogenetic trees and ancestral host reconstructions.** Alpha-CoV (A) and beta-CoV (B) maximum  
783 clade credibility annotated trees using complete datasets of *RdRp* sequences and bat host family as  
784 discrete character state. Pie charts located at the root and close to the deepest nodes show the state  
785 posterior probabilities for each bat family. Branch colors correspond to the inferred ancestral family  
786 with the highest probability. Branch lengths are scaled according to relative time units (clock rate = 1.0).  
787 Well-supported nodes (posterior probability > 0.95) are indicated with a black dot. The ICTV approved  
788 CoV subgenera were highlighted: *Rhinacovirus* (L1), *Decacovirus* (L2), *Myotacovirus* (L3), *Pedacovirus*  
789 (L5), *Nyctacovirus* (L6), *Minunacovirus* (L7) and an unidentified lineage (L4) for alpha-CoVs; and  
790 *Merbecovirus* (Lineage C), *Nobecovirus* (lineage D), *Hibecovirus* (lineage E) and *Sarbecovirus* (Lineage B)  
791 for beta-CoVs.

792 **Fig. 3 Phylogenetic relationships within the *Sarbecovirus* subgenus (beta-CoVs).** Maximum clade  
793 credibility tree (A) including 202 *RdRp* sequences from the *Sarbecovirus* subgenus isolated in bats, two  
794 sequences of SARS-CoV-2 and one sequence of SARS-CoV isolated in humans and eight sequences  
795 isolated in Malayan pangolins (*Manis javanica*). Well-supported nodes (posterior probability > 0.95) are  
796 indicated with a black dot. Tip colors correspond to the host genus, SARS-CoV-2 sequences and SARS-  
797 CoV sequence are highlighted in grey and black, respectively. Median-joining network (B) including 202  
798 *RdRp* sequences from the *Sarbecovirus* lineage isolated in bats, two sequences of SARS-CoV-2 and one

799 sequence of SARS-CoV isolated in humans and eight sequences isolated in Malayan pangolins (*Manis*  
800 *javanica*). Colored circles correspond to distinct CoV sequences, circle size is proportional to the number  
801 of identical sequences in the data set. Small black circles represent median vectors (ancestral or  
802 unsampled intermediate sequences). Branch length is proportional to the number of mutational steps  
803 between haplotypes.

804 **Fig. 4 Inter-family host switches.** Strongly supported host switches between bat families for alpha- (A)  
805 and beta-CoVs (B). Arrows indicate the direction of the switch; arrow thickness is proportional to the  
806 switch significance level, only host switches supported by strong Bayes factor (BF) > 10 are shown.  
807 Histograms of total number of host switching events (state changes counts using Markov jumps) from/to  
808 each bat family along the significant inter-family switches for alpha- (C) and beta-CoVs (D).

809 **Fig. 5 Inter-genus host switches.** Strongly supported host switches between bat genera for alpha- (A)  
810 and beta-CoVs (B) and their significance level (Bayes factor, BF). Only host switches supported by strong  
811 BF values > 10 are shown. Line thickness is proportional to the switch significance level. Red lines  
812 correspond to host switches among bat genera belonging to different families, black lines correspond to  
813 host switches among bat genera from the same family. Arrows indicate the direction of the switch.  
814 Genus names are colored according to the family they belong to using the same colors as in Fig. 2 and 3.

815 **Fig. 6 CoV spatiotemporal dispersal in China.** Strongly supported dispersal routes (Bayes factor, BF > 10)  
816 over recent evolutionary history among China zoogeographic regions for alpha- (A) and beta-CoVs (B).  
817 Arrows indicate the direction of the dispersal route; arrow thickness is proportional to the dispersal  
818 route significance level. Darker arrow colors indicate older dispersal events. Histograms of total number  
819 of dispersal events (Markov jumps) from/to each region along the significant dispersal routes for alpha-  
820 (C) and beta-CoVs (D). NO, Northern region; CN, Central northern region; SW, South western region; CE,  
821 Central region; SO, Southern region; HI, Hainan island.

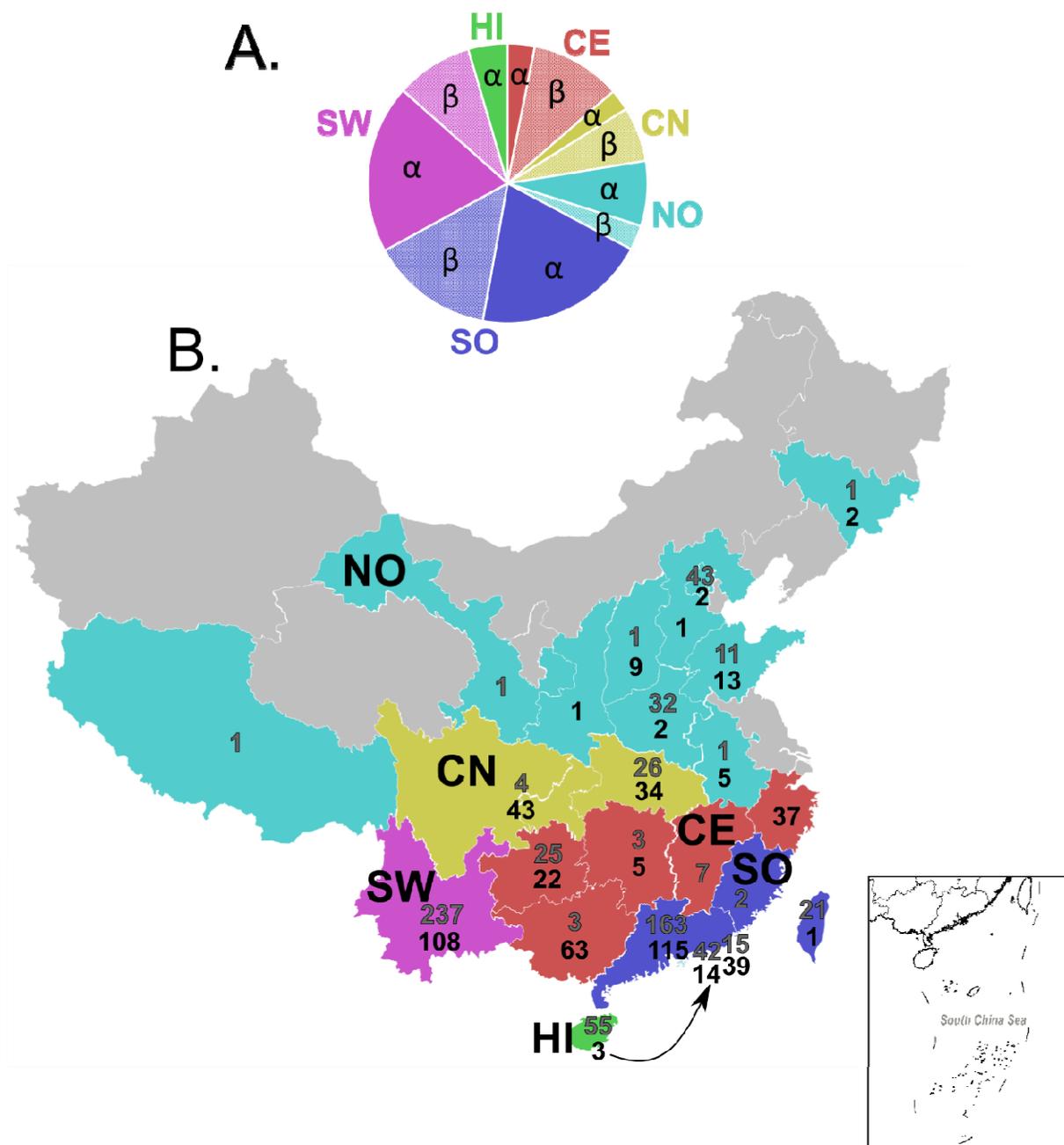
822 **Fig. 7 Phylogenetic diversity.** Metrics of CoV phylogenetic diversity within each bat family (A), genus (B)  
823 and zoogeographic regions (C): standardized effect size of Mean Phylogenetic Distance (SES MPD), on  
824 the left panels; and standardized effect size of Mean Nearest Taxon Distance (SES MNTD), on the right  
825 panels. One-tailed p-values (quantiles) were calculated after randomly reshuffling tip labels 1000 times  
826 along the entire phylogeny. Values departing significantly from the null model (p-value < 0.05) are  
827 indicated with an asterisk, all exact p-values are available in Supplementary Tables 14-27. NO, Northern  
828 region; CN, Central northern region; SW, South western region; CE, Central region; SO, Southern region;  
829 HI, Hainan island.

830 **Fig. 8 Phylogenetic diversity.** Standardized effect size of Mean Phylogenetic Distance (SES MPD) and  
831 phylogenetic ordination among bat host families (A, B) and genera (C, D) for alpha- and beta-CoVs.  
832 Boxplots for each host family and genus show the mean (cross), median (dark line within the box),  
833 interquartile range (box), 95% confidence interval (whisker bars), and outliers (dots), calculated from all  
834 pairwise comparisons between bat families (n=10 for alpha-CoVs and n=6 for beta-CoVs) and genera  
835 (n=91 for alpha-CoVs and n=105 for beta-CoVs).

836 **Fig. 9 Phylogenetic diversity.** Standardized effect size of Mean Phylogenetic Distance, SES MPD) and  
837 phylogenetic ordination among zoogeographic regions for alpha- (A) and beta-CoVs (B). Boxplots for  
838 each region show the mean (cross), median (dark line within the box), interquartile range (box), 95%  
839 confidence interval (whisker bars), and outliers (dots), calculated from all pairwise comparisons between  
840 regions (n=15 for alpha-CoVs and n=10 for beta-CoVs). NO, Northern region; CN, Central northern  
841 region; SW, South western region; CE, Central region; SO, Southern region; HI, Hainan island.

842

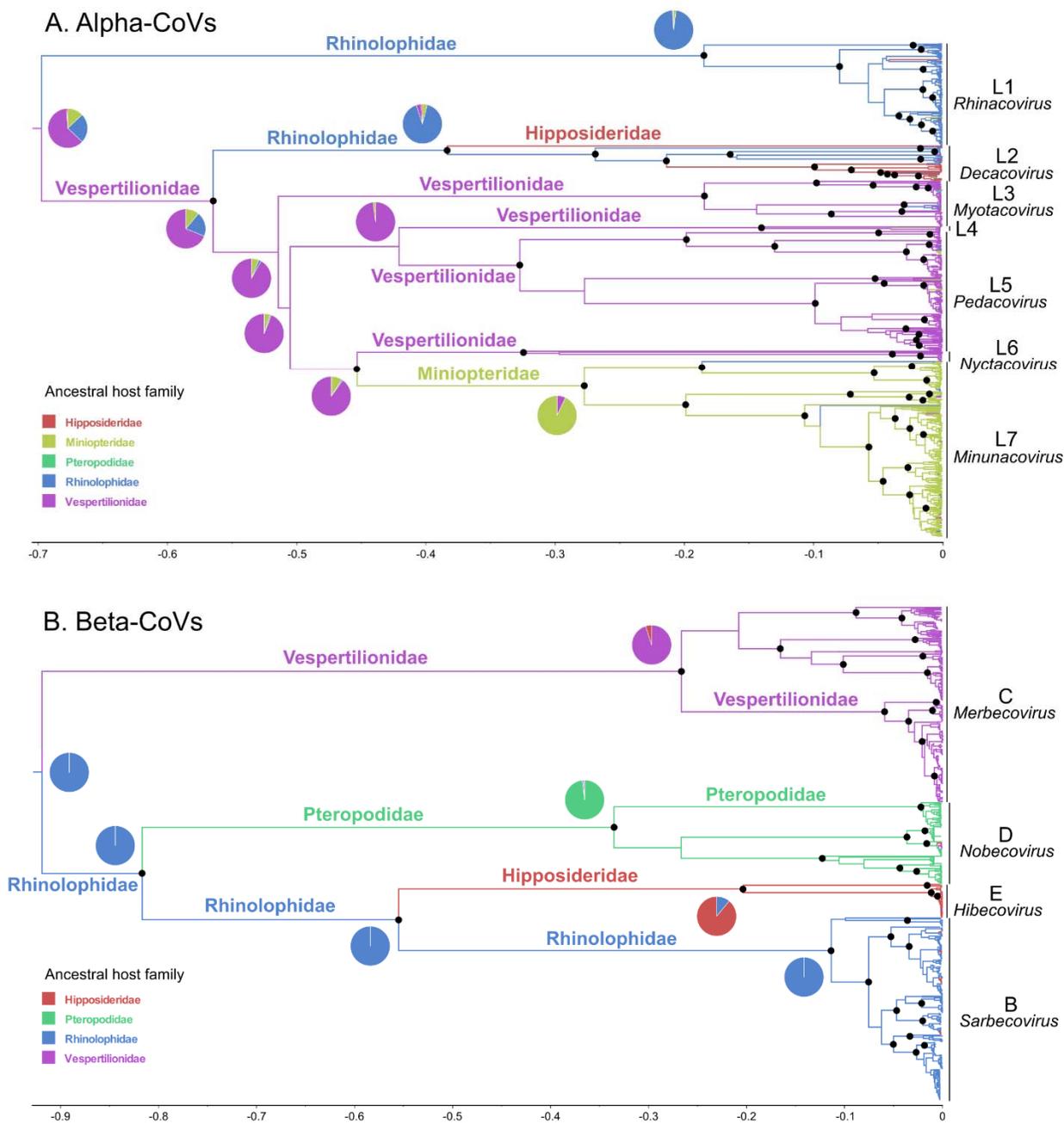
843 **Figure 1**



844

845

846 **Figure 2**



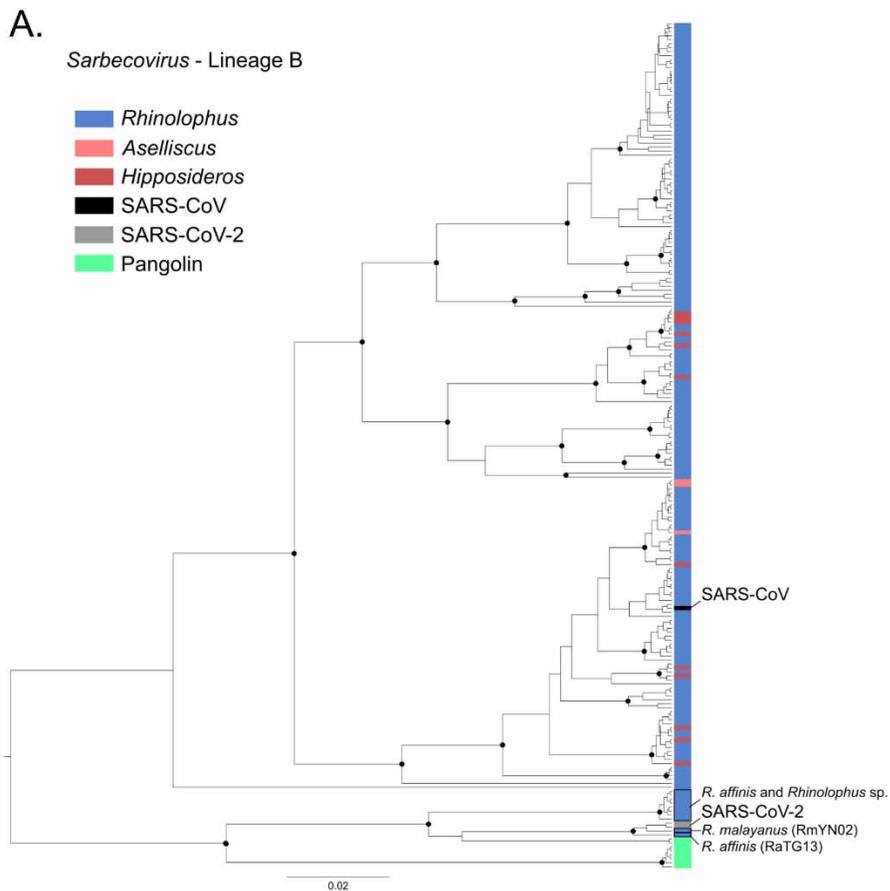
847

848

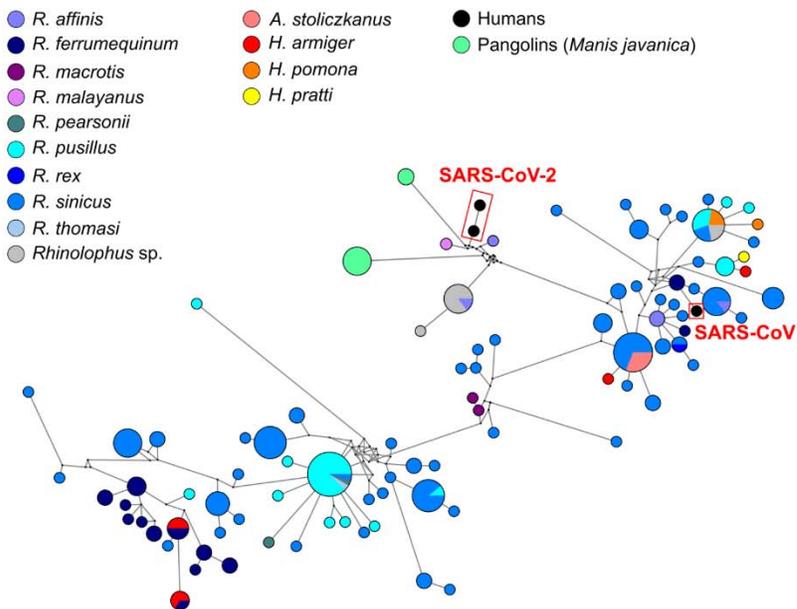
849

850

851 **Figure 3**

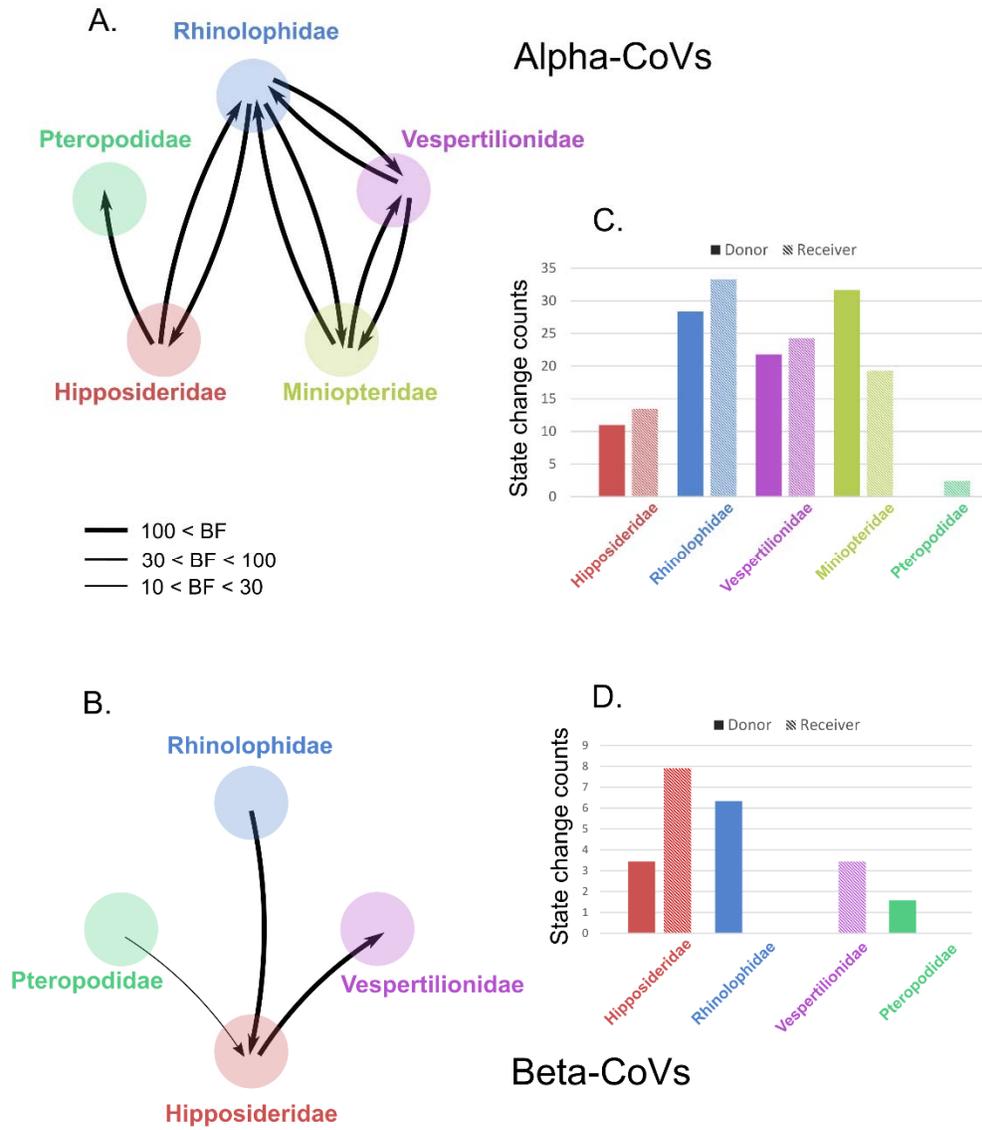


**B.**



852

853 **Figure 4**

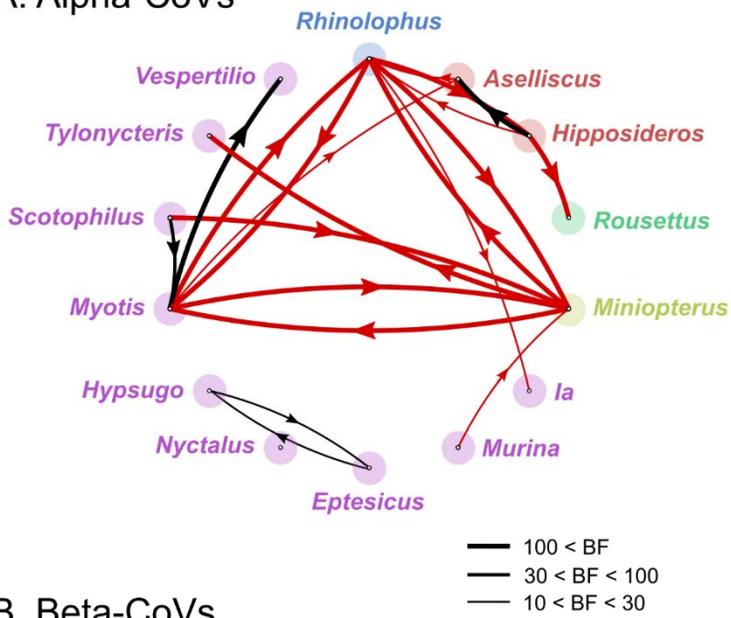


854

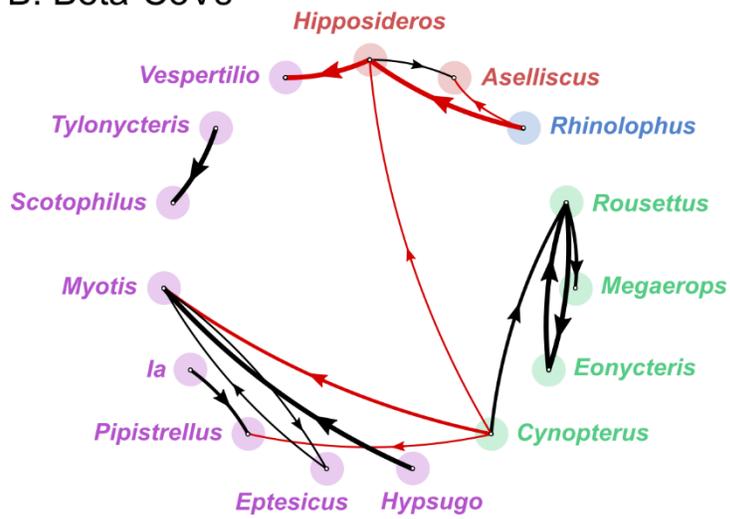
855

856 **Figure 5**

**A. Alpha-CoVs**



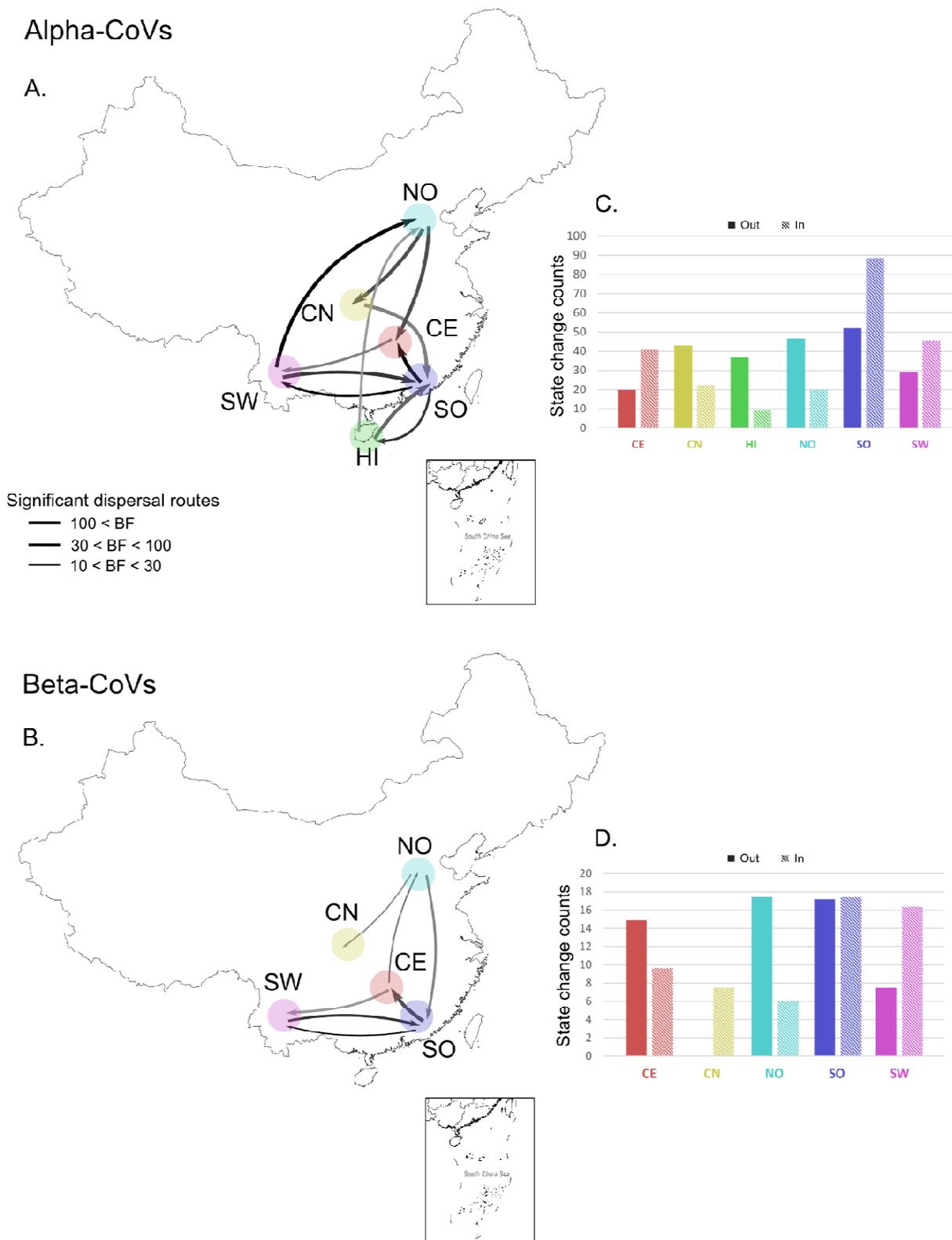
**B. Beta-CoVs**



857

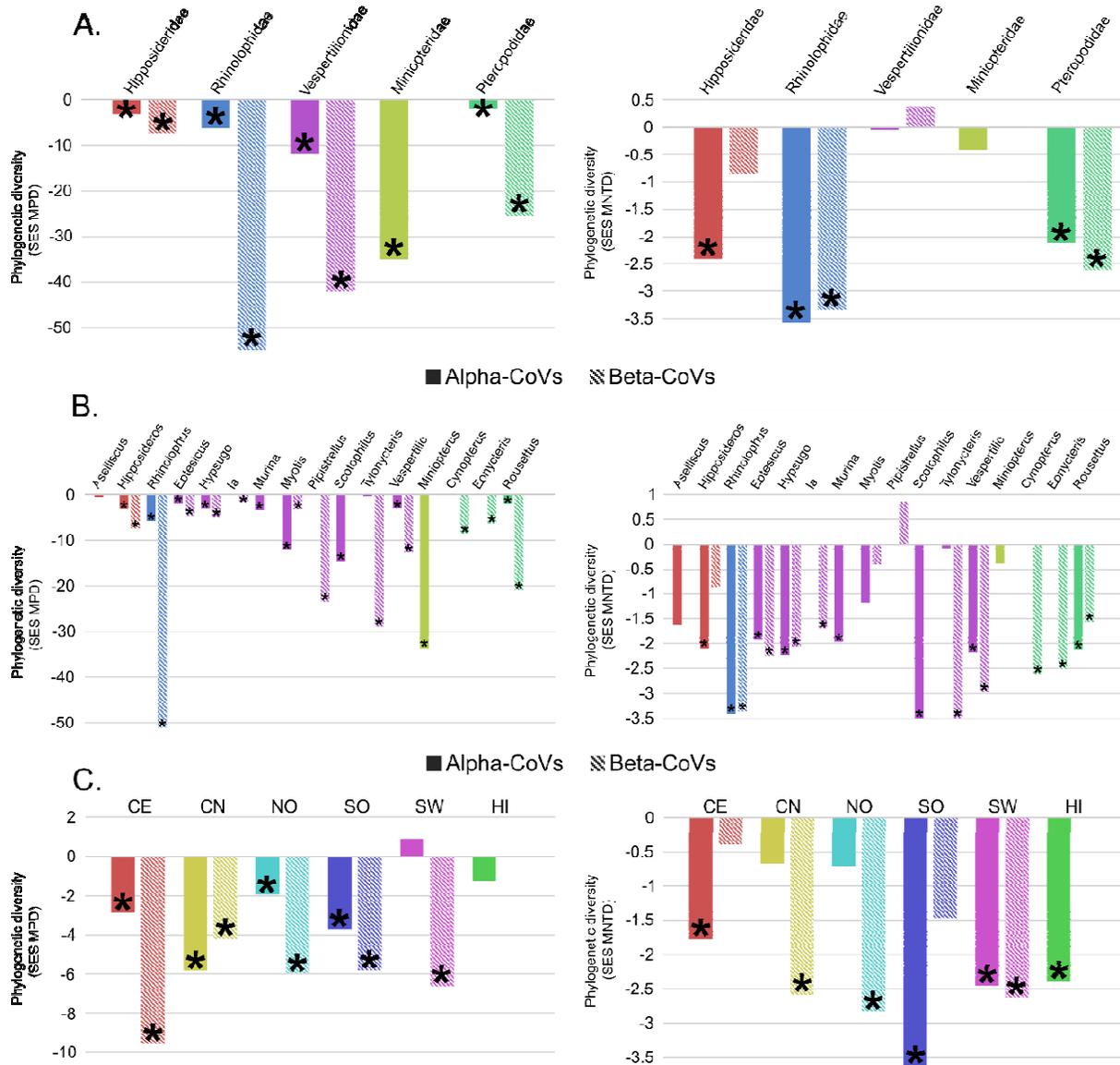
858

859 **Figure 6**



860

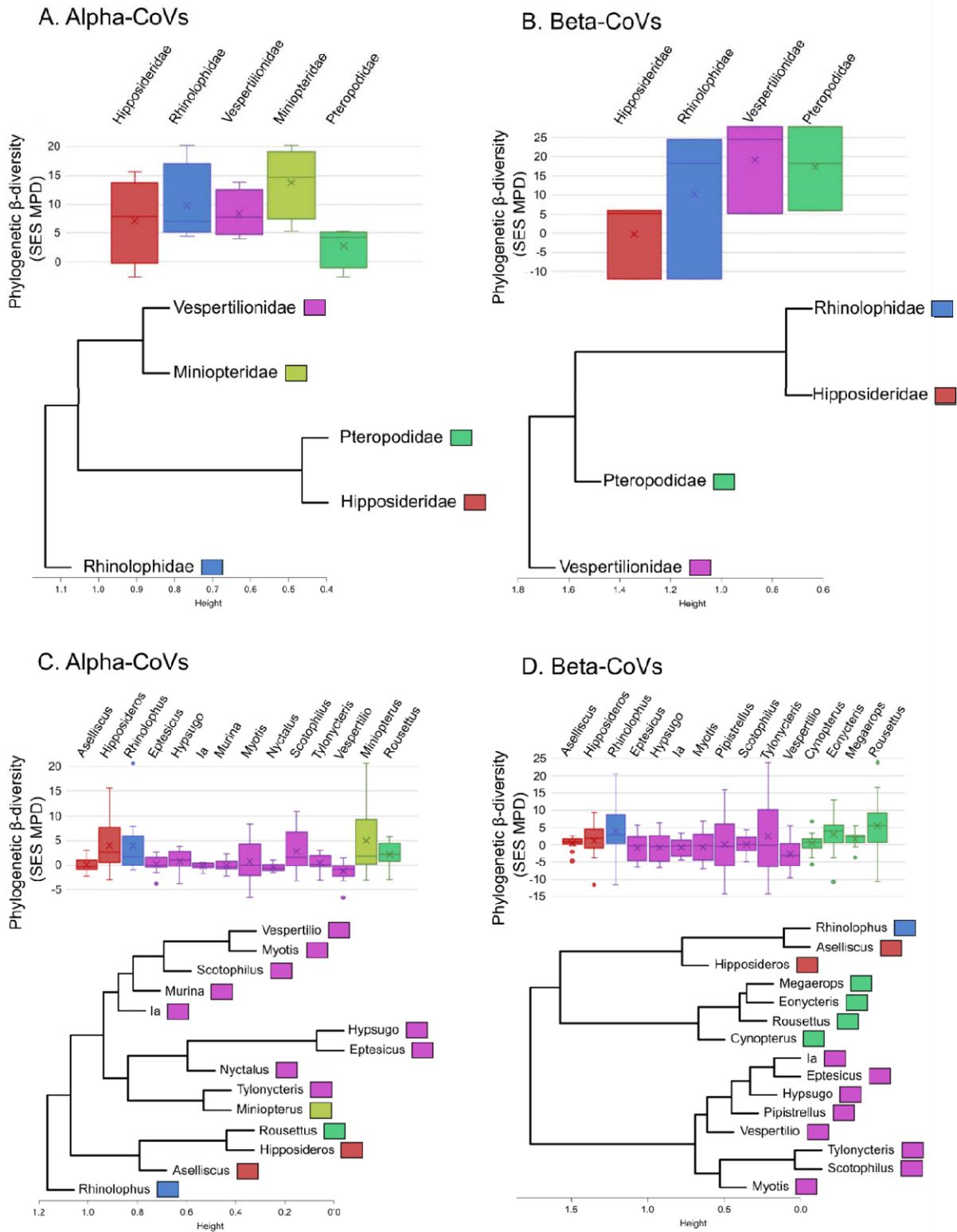
861 **Figure 7**



862

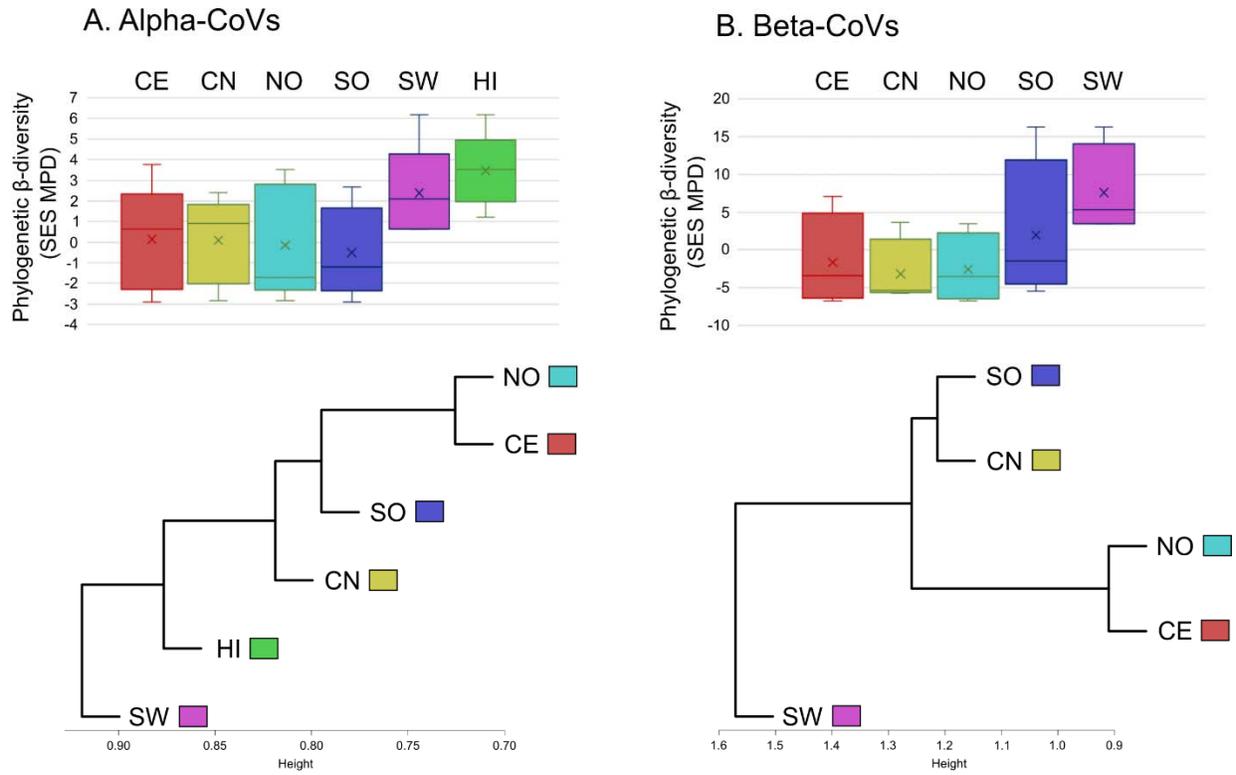
863

864 **Figure 8**



865

866 **Figure 9**



867

868