

1                   **Genome-wide data inferring the evolution and population**  
2                   **demography of the novel pneumonia coronavirus (SARS-CoV-2)**

3

4   Bin Fang<sup>1\*</sup>, Linlin Liu<sup>1\*</sup>, Xiao Yu<sup>1\*</sup>, Xiang Li<sup>1\*</sup>, Guojun Ye<sup>1</sup>, Juan Xu<sup>3</sup>, Ling Zhang<sup>3</sup>,  
5                   Faxian Zhan<sup>1</sup>, Guiming Liu<sup>4</sup>, Tao Pan<sup>2#</sup>, Yilin Shu<sup>2#</sup>, Yongzhong Jiang<sup>1#</sup>

6

7                   1. Influenza Reference Laboratory, Institute of Health Inspection and Testing, Hubei  
8                   Provincial Center for Disease Control and Prevention, Wuhan 430079, China

9                   2. Anhui Province Key Laboratory for Conservation and Exploitation of Biological  
10                   Resource, College of Life Sciences, Anhui Normal University, Wuhu 241000, Anhui,  
11                   China

12                   3. Hubei Province Key Laboratory of Occupational Hazard Identification and Control,  
13                   School of Public Health, Wuhan University of Science and Technology, Wuhan,  
14                   Hubei, 430065, China

15                   4. Beijing Agro-Biotechnology Research Center, Beijing Academy of Agriculture and  
16                   Forestry Sciences, 100097, Beijing, China

17                   \*These authors have equal contribution to this study.

18                   #Correspondence: Tao Pan, [pantao20070778@163.com](mailto:pantao20070778@163.com); Yilin Shu, [sjshuyilin@163.com](mailto:sjshuyilin@163.com);

19                   Yongzhong Jiang, [hbcddxd@163.com](mailto:hbcddxd@163.com).

20

21

22

23 **Other Author's Email:**

24 Bin Fang: nicolfang@163.com

25 Linlin Liu: mice0809@163.com

26 Xiao Yu: fish-yuxiao@hotmail.com

27 Xiang Li: li\_xiang75@yeah.net

28 Guojun Ye: gjy0822@sina.com

29 Juan Xu: xujuan9270@163.com

30 Ling Zhang: zhangling@wust.edu.cn

31 Faxian Zhan: [zhanfx@163.com](mailto:zhanfx@163.com)

32 Guiming Liu: mingguiliu@aliyun.com

33

34

35 **Abstract**

36 Since December 2019, coronavirus disease 2019 (COVID-19) emerged in  
37 Wuhan, Central China and rapidly spread throughout China. Up to March 3, 2020,  
38 SARS-CoV-2 has infected more than 89,000 people in China and other 66 countries  
39 across six continents. In this study, we used 10 new sequenced genomes of  
40 SARS-CoV-2 and combined 136 genomes from GISAID database to investigate the  
41 genetic variation and population demography through different analysis approaches  
42 (e.g. Network, EBSP, Mismatch, and neutrality tests). The results showed that 80  
43 haplotypes had 183 substitution sites, including 27 parsimony-informative and 156  
44 singletons. Sliding window analyses of genetic diversity suggested a certain mutations

45 abundance in the genomes of SARS-CoV-2, which may be explaining the existing  
46 widespread and high adaptation of the deadly virus. Phylogenetic analysis showed  
47 that the view, pangolin acted as an intermediate host, may be controversial. The  
48 network indicated that, in the original haplotype (H14), one patient sample lived  
49 near the Huanan seafood market (approximate 2 km), indicating high possibility of  
50 the patient having a history of unconscious contact with this market. However, based  
51 on this clue, we cannot accurately concluded that whether this market was the origin  
52 center of SARS-CoV-2. Additionally, 16 genomes, collected from this market,  
53 assigned to 10 haplotypes, indicated a circulating infection within the market in a  
54 short term and then leading to the outbreak of SARS-CoV-2 in Wuhan and other areas.  
55 The EBSP results showed that the first estimated expansion date of SARS-CoV-2  
56 began from 7 December 2019, which may indicated that the transmission could have  
57 begun from person to person in mid to late November.

58

59 **Key words** SARS-CoV-2; metagenomic next-generation sequencing; virus  
60 evolution; population demography; phylogenetic relationship

61

## 62 **Introduction**

63 As the largest non-segmented genomes among all the RNA viruses (about 30 kb  
64 in length), Coronaviruses (CoVs) own the plasticity due to the mutation and  
65 recombination, which increased the potential risks of spread across species [1, 2]. The  
66 COVID-19 (original named 2019-nCoV) is the seventh member of enveloped RNA

67 coronavirus (subgenus, Sarbecovirus; subfamily, Orthocoronavirinae) [3]. In early  
68 December, 2019, an unexplained pneumonia associated with the severe acute  
69 respiratory syndrome coronavirus 2 (SARS-CoV-2) emerged in Wuhan, China [4, 5].  
70 The rapid emergence and spread of the SARS-CoV-2 between infected and healthy  
71 people became so devastating as a large population within Wuhan was getting  
72 infected [6, 7]. There's evidence that SARS-CoV-2 may have originated from bats, but  
73 there is no clear information about the intermediate host that transferred it to humans  
74 [8-10]. On 30 January 2020, the World Health Organization (WHO) declared the  
75 outbreak of COVID-19 to be a Public Health Emergency of International Concern. Up  
76 to March 3, 2020, SARS-CoV-2 has infected more than 89,000 people in China and  
77 other 66 countries across the six continents (source: World Health Organization  
78 report).

79 Despite the worldwide rapid spread, the genomic variation dynamics,  
80 evolutionary rate, and virus transmission dynamics of SARS-CoV-2 are not yet well  
81 understood. In several recent studies, phylogenetic relationships, variations,  
82 evolutionary rates, and propagation dynamics were analyzed using limited genomic  
83 data from the SARS-CoV-2 [5, 9-17]. As the epidemic progresses, many research  
84 institutes around the world have obtained and submitted SARS-CoV-2 genome  
85 sequences, increasing the number of SARS-CoV-2 genome sequences in the Data  
86 Centers (i.e, GISAID). Given the extremely rapid spread of SARS-CoV-2, an updated  
87 analysis with significantly larger sample sizes by incorporating cases throughout the  
88 world is urgently needed. This would allow for more accurate identification of

89 SARS-CoV-2 variant dynamics, evolutionary rates, virus transmission dynamics, and  
90 epidemic history in order to effectively implement public health measures, promote  
91 drugs and vaccines development and to efficiently prevent similar epidemics in the  
92 future.

93       Recent research findings have suggested bat and other non-bat intermediate  
94 mammals (such as pangolin) as potential intermediate host for the SARS-CoV-2 that  
95 have so far been widely transmitted to humans [5, 9, 18]. According to the medical  
96 information of the first patients to be infected in Wuhan, 27 out of 41 patients were  
97 found to be traders of wild mammals in the Huanan seafood market, which directly  
98 suggests that the Huanan seafood market as the possible source of origin of the  
99 SARS-CoV-2, and later got transmitted it to other areas by the infected people [4, 19].  
100 However, due to some infected persons having no relation or access to the Huanan  
101 food market, some researchers are skeptical of the market as the only actual origin or  
102 the source of SARS-CoV-2 transmission to humans [4, 19, 20]. Therefore, up to date,  
103 in the absence of key potential intermediary host, the origin and transmission pattern  
104 of SARS-CoV-2 still need to be thoroughly probed for a more reliable and accurate  
105 information on the origin and transmission mechanisms of the virus.

106       In this study, 56 full genome sequences of outgroups and 136 genomes of  
107 SARS-CoV-2 from GISAID EpiFluTM database (access date 22 February 2020) were  
108 collected. Ten new sequenced genomes of SARS-CoV-2 were obtained by  
109 macrogenomic sequencing. The time origin, genetic diversity, transmission dynamics  
110 and evolutionary history of SARS-CoV-2 were analyzed based on the genomic data to

111 provide understanding on the origin, transmission pathway, and evolutionary  
112 characteristics of SARS-CoV-2 outbreak.

## 113 **Materials and methods**

### 114 *Patients and samples*

115 In this study, ten patients with unexplained viral pneumonia from five hospitals  
116 in Hubei Province were included. Four samples were collected in December, four in  
117 January as well as two samples in February. The December samples were examined  
118 by LightCycle 480II fluorescent PCR instrument (Roche, Basel, Switzerland) for  
119 investigating 26 respiratory pathogens and detection of the virus using SARS-Cov and  
120 MERS-Cov primer probes. All the ten samples were detected by the SARS-CoV-2  
121 primer probes. Two samples were used in the virus isolation. The 10 samples were  
122 also subjected to metagenomic next-generation sequencing. Specific sample  
123 information is presented in Table 1.

### 124 *Virus isolation*

125 Virus were isolated by HUH7 cells from one lavage fluid sample and one  
126 pharynx swab sample. The cells were monitored daily for cytopathic effects by light  
127 microscopy. HUH7 cells were harvested after 6 days of culturing at 37°C. The  
128 supernatant was collected to extract the total RNA nucleic acid through EZ1 virus  
129 mini kit v.2.0 (955134; Qiagen, Heiden, Germany), detection was performed by Light  
130 Cycle 480II fluorescent PCR instrument (Roche, Basel, Switzerland) followed by  
131 subsequent sequencing.

### 132 *Library preparation and sequencing*

133 Total RNA extracted from 10 samples were subjected to metagenomic next  
134 generation sequencing testing by EZ1 virus mini kit v.2.0 (955134; Qiagen, Heiden,  
135 Germany). TruSeq Stranded Total RNA Library Prep Kit (Illumina, San Diego, CA,  
136 USA) was used to remove rRNA, reverse transcription and synthesize the  
137 double-stranded DNA. The fragmenting, modifying the ends, connecting the joints  
138 and enriching of DNA were done by NexteraXT library prepkit (Illumina, San Diego,  
139 CA, USA), then the DNA library was obtained. Iseq<sup>TM</sup> 100 i1 Cartridge, Miseq v2  
140 reagent kit, and High output Reagent Cartridge (illumina, San Diego, CA, USA) were  
141 used for deep sequencing in Iseq100, Miseq, and Miniseq platforms (illumina, San  
142 Diego, CA, USA), respectively. About 11.5 GB data were obtained for each sample.

#### 143 *Virus genome analysis*

144 Using CLC Genomics Workbench 12 and Geneious 12.0.1 software (QIAGEN  
145 Bioinformatics, Redwood City, CA, USA) and using reference sequence  
146 BetaCoV/Wuhan-Hu-2019 (EPI\_ISL\_402125), the raw sequencing data of Illumina  
147 were analyzed. The open reading frames of the verified genome sequences were  
148 predicted using Geneious and annotated using the Conserved Domain Database [21].

#### 149 *Phylogenetic reconstructions*

150 In this study, for probing the evolutionary history of SARS-CoV-2 (Table S1),  
151 136 complete genomes from GISAID (Table S1, access on 22 February, 2020) and  
152 our new sequenced 10 genomes of this virus were included. In total, 56 outgroups  
153 were collected following the previous study [5]. Mafft v.7.450 was used to align the  
154 sequence of SARS-CoV-2 with reference sequences [22]. Phylogenetic analyses of  
155 the genomes were done with RAxML v.8.2.9 [23] in 1000 bootstrap replicates,

156 employing the general time reversible nucleotide substitution model. The tree was  
157 visualized with FigTree v.1.4.3 [24].

### 158 ***Population structure***

159 On this analysis technique, to precisely decode the evolutionary history of  
160 SARS-CoV-2, the genome EPI ISL 402131 (bat-RaTG13-CoV) from GISAID was  
161 also included as the outgroup following the previous study [25], because it is the  
162 closest sister betacoronavirus to SARS-CoV-2.3. The 136 complete genome  
163 sequences of SARS-CoV-2 and an outgroup (bat-RaTG13-CoV) were aligned using  
164 MAFFT then the alignment was manually checked using Geneious. For probing the  
165 haplotypes relationships among localities, the Network v.4.6.1 (Fluxus Technology,  
166 Suffolk, UK) was used to construct the minimum-spanning network based on the full  
167 median-joining algorithm [26]. During the alignment, 10 genomes containing  
168 ambiguous sites or “N” bases or more degenerated bases were excluded in this study  
169 (Table S1). In addition, 4 genome sequences (EPI ISL 409067, EPI ISL 412981, EPI  
170 ISL 407071 and EPI ISL 408489) with one degenerated base were included, which  
171 were divided into two sequences without degenerated base. For ensuring the accuracy  
172 of this analysis, the 5’UTR and 3’UTR contain missing and ambiguous sites of both  
173 regions were excluded in the following alignment analyses. DnaSP v.5.10 [27] was  
174 used to convert the relevant format. Population genetic indices in was estimated in  
175 DnaSP, including nucleotide diversity ( $\pi$ ) and haplotype diversity ( $Hd$ ).  $F_{ST}$  between  
176 haplotypes was calculated in Mega 6.0 based on haplotype frequency differences with  
177 10,000 permutations [28]. Additionally, a sliding window of 500 bp and steps of 50  
178 bp were used to show nucleotide diversity ( $\pi$ ) for the entire alignment this data.  
179 Nucleotide diversity for the entire alignment was plotted against midpoint positions of



180 each window. To indicate the relative position of the mutations in the genome, we  
181 selected the EPI ISL 402124 sequence as the reference genome.

## 182 *Population demography*

183 The alignment was then imported into DnaSP for haplotype analyses. Population  
184 size changes were estimated based on a constant population size hypothesis using  
185 DnaSP, in combination with neutrality tests (Tajima's *D* and *Fu*'s *F<sub>s</sub>*). The mismatch  
186 distribution was estimated based on the genomes of SARS-CoV-2 in Arlequin to test  
187 the hypothesis of recent population growth. The Harpending's raggedness index (*Rag*)  
188 and the sum of squares deviation (*SSD*) were used to determine the smoothness of  
189 observed mismatch distribution and the degree of fit between observed and simulated  
190 data [29]. Because of the sensitivity to demographic changes of neutral tests, Tajima's  
191 *D* [30] and *Fu*'s *F<sub>s</sub>* [31] were estimated using 10,000 coalescent simulations to assess  
192 the significance in Arlequin. In addition, the Extended Bayesian skyline plot (EBSP)  
193 analysis was conducted to examine past population dynamics of SARS-CoV-2 based  
194 on 146 genomes by BEAST v.1.8 [32]. The divergence times were estimated in  
195 BEAST using a Bayesian Markov chain Monte Carlo (MCMC) method with a strict  
196 clock. In this study, the substitution rate was set as  $0.92 \times 10^{-3}$  (95% CI,  
197  $0.33 \times 10^{-3}$ - $1.46 \times 10^{-3}$ ) substitution/site/year based on the most recent estimation for  
198 SARS-CoV-2 [25]. The other parameters were set as follows: extended Bayesian  
199 skyline process, 10 million MCMC generations, sampling every 1,000<sup>th</sup> iteration, the  
200 initial 25% burn-in. Tracer was used to check the convergence of the MCMC analyses  
201 (effective sample size [ESS] values >200). Convergence of the two independent  
202 MCMC runs was assessed in Tracer, as was convergence of model parameter values  
203 (ESS) to ensure ESS values >200.

## 204 **Results and discussion**

### 205 *Genomic variations of SARS-CoV-2*

206 According to the database, the genome size of SARS-CoV-2 varied from 29,409  
207 bp to 29,911 bp. In the network dataset, the aligned matrix was 29,130 bp in length  
208 including 183 variable sites, of which 27 were parsimony-informative and 156 were  
209 singletons, which were classified as 80 haplotypes (Table S2). Nucleotide diversity ( $\pi$ )  
210 was  $0.16 \times 10^{-3} \pm 0.02 \times 10^{-3}$  (Standard Deviation, SD) and haplotype diversity ( $Hd$ ) was  
211  $0.954 \pm 0.013$  (SD) and variance of  $Hd$  was  $0.16 \times 10^{-3}$  (Table S3). In this study, our  
212 new sequenced 10 genomes of SARS-CoV-2 represented 10 haplotypes (H3, H8, H10,  
213 H14, H16, H28, H29, H30, H76 and H78; Table S1). Among these 10 genomes, four  
214 samples were collected from those patients who had been exposed to the Huanan  
215 seafood market in Wuhan. Additionally, the HBCDC-HB-03/2019 belonged to the  
216 core H3 haplotype and most of them represented new haplotype. Sliding window  
217 analysis of SARS-CoV-2 revealed significant regional variation across the alignment  
218 (Fig. 1). The plot readily showed a relative high degree of nucleotide variation  
219 amongst the aligned SARS-CoV-2 genomes for any given window of 500 bp and  
220 steps of 50 bp, with the  $\pi$  ranging from 0.00005 to 0.0014. Additionally, based on the  
221 curve, there were several relative high variation area, compared with other fragments  
222 (Fig. 1). The  $F_{ST}$  analyses among 80 haplotypes ranged from 0.00003 (e.g, H3 and H9;  
223 H14 and H15; H3 and H20) to 0.00134 (H19 and H80), which indicated genetic  
224 differentiation among these haplotypes (Fig. 2, Table S4). Overall, the H19 showed  
225 relatively large genetic distance with other haplotypes, while the H3 showed opposite  
226 pattern, and was also confirmed by the network figure (Fig. 2, Table S4). Overall,  
227 SARS-CoV-2 maintain the relatively rich mutations that has occurred across the world,  
228 which may be the main reason of the numerous subgenomic RNAs generated during

229 the viral replication [3, 9]. This phenomenon also provides the possibility of  
230 widespread and adaptation for this virus. Fortunately, under the strict quarantine  
231 policy in China since 23 January 2020, the circulation and spreading of some  
232 haplotypes may have relatively reduced in circulation frequency relatively, compared  
233 with the early stage of COVID-19.

#### 234 ***Phylogenetic relationship of SARS-CoV-2***

235 The virus's popularity and its intermediate host has been a topic of great concern.  
236 Some researchers have suggested that the SARS-CoV-2 originated from the bats [5,9].  
237 However, some other researchers have proposed that the non-bat intermediate  
238 mammals (e.g. pangolins) may be the transmission path of this genus [18]. However,  
239 other studies have also failed to confirm the pangolins as the intermediate host to this  
240 viral genus [33]. In this study, based on the phylogenetic tree, bat-RaTG13-CoV was  
241 the sister group to the SARS-CoV-2 and the two pangolin samples collected from  
242 Guangdong which showed a relatively large genetic distance from the SARS-CoV-2  
243 (Fig. 3 and S1). Based on this finding, , the study suggested that pangolin may not be  
244 the intermediate host, as had been reiterated by other previous study [33]. As the  
245 Huanan seafood market was closed on 1 January 2020, this created many challenges  
246 in identifying the first intermediate hosts (whether people or animal) of SARS-CoV-2.  
247 Therefore, more sampling and analysis of the sample from this area may suggest  
248 clearer results for more concrete conclusions in future studies.

#### 249 ***Evolutionary relationships of SARS-CoV-2 haplotypes***

250 The evolutionary network of 80 haplotypes of SARS-CoV-2, with bat-RaTG13-CoV  
251 as the outgroup, is shown in Figure 4. The network analysis showed typical star  
252 network, with several core haplotype node (H3, H14, H15) and edge haplotype nodes.  
253 In the network, fifty-five satellite haplotypes and H14 connected to the H3 haplotype;

254 eighteen satellite haplotypes and H15+H3 connected to the H14 haplotype; and five  
255 satellite haplotypes and H14 connected to H15. Most of these haplogroups were  
256 separated by one to six mutations, except two haplotypes (H19 and H80, Fig. 4). The  
257 evolutionary network showed that bat-RaTG13-CoV to be connected through a  
258 hypothesized haplotype (mv6) to the H15 and H31 haplotypes by single mutations  
259 (Fig. 4). However, the mutations between mv6 and bat-RaTG13-CoV was more than  
260 1,000, which indicated that SARS-CoV-2 still has a relative distant kinship with the  
261 outgroup (bat-RaTG13-CoV).

262 As the most abundant haplotype, H3 included 28 samples, while 55 satellite  
263 haplotypes are directly derived from the H3 haplotype (Fig. 4 and Table S1).  
264 Moreover, many haplotypes from other countries should also be derived from the H3  
265 haplotype. The current samplings showed that the H3 haplotype has been found to be  
266 in 28 samples, but 12 of 28 samples were collected from Wuhan in Hubei Province.  
267 Fifteen of the 55 satellite haplotypes were also collected from Hubei Province (Fig. 4).  
268 One possible explanation was that a common haplotype from the Huanan seafood  
269 market (Fig. 4) was rapidly circulated at an early stage of human-to-human  
270 transmissions. It is worth noting that in the H14, there are two samples (EPI ISL  
271 406801 and EPI ISL 412979) from Wuhan. Although these two hosts didn't directly  
272 link with the Huanan market, one of the host (EPI ISL 412979) having lived in a  
273 residential area about 2 kilometers from the Huanan seafood market. This also  
274 provides the possibility of them having indirect possible contact with the Huanan  
275 seafood market.

276 Additionally, the networks of 24<sup>th</sup> December, 2019 - 30<sup>th</sup> December, 2019 and  
277 24<sup>th</sup> December, 2019 - 6<sup>th</sup> January, 2020 also suggested that the most frequent  
278 haplotype (H3) occurred. All the haplotypes were collected from Wuhan, which

279 indicated that Wuhan may have acted as the important origin center (Fig. 4). On the  
280 other hand, in this study, 16 sequences of Huanan seafood market (4 new sequenced  
281 and 12 GISAID data) were collected. Fifteen of out of the sixteen samples were  
282 collected before 1<sup>st</sup> Januray, 2020 and one samples was collected on 8<sup>th</sup> Januray, 2020.  
283 These samples were represented ten haplotypes (H2, H3, H4, H5, H6, H7, H8, H10,  
284 H11, H16), which indicated the high proportion of haplotype diversity in Huanan  
285 seafood market. Among those haplotypes, the H3 haplotype, the most abundant, was  
286 present in 6 samples from Huanan seafood market, while the other haplotypes were  
287 directly derived from the H3 haplotype. All the samples from the Huanan seafood  
288 market had the H3 haplotype and other 9 derived haplotypes (Fig. 4), indicating that  
289 there were circulated infections within the market in a short term. Noteworthy, in the  
290 network, a total of 65 virus samples from 15 other countries were assigned to 43  
291 haplotypes. Among them, 27 haplotypes were satellite haplotypes of H3 haplotype, 11  
292 haplotypes were satellite haplotypes of H14 haplotype, and 3 haplotypes were satellite  
293 haplotype of H15 haplotype. Indeed, most of haplotypes originated from the H3,  
294 which may hint the phenomenon was related to the input of virus carriers from  
295 Wuhan.

296       Based on the outgroup (bat-RaTG13-CoV) having a possible direct connection  
297 with outgroup H15 and H31, it indicates that both the H15 and H31 were the  
298 suggested ancestral haplotypes (Fig. 4), as had also been suggested by other previous  
299 study [25]. The H15 was only recovered from five Shenzhen (Guangdong) samples,  
300 while H31 included 3 United States samples and 1 Fujian sample from China (Fig. 4).  
301 Based on the epidemiological statistics, all the patients from the H15 were ever  
302 traveled to Wuhan [16]. As to the H31, three genomes from the same patient in Unite  
303 State [34] and 1 genome from Fujian of China (Table S1) were included. For the

304 patient in United State, He may have infected during the period of visiting his family in  
305 China [25]. The travel history of Fujian sample was unclear, but we cannot deny that  
306 it has a history of unconscious Wuhan contacts. Based on the theory in the previous  
307 study [25], this current study also supported two main evolutionary paths, that the  
308 available haplotypes could be from H15 through H14 to H3, or from H31 through  
309 H14 to H3 (Fig. 4). Both these two paths demonstrate that H14 was the key connection  
310 from an ancestral haplotype to H3. However, unlike the previous research, in H14,  
311 two individuals (EPI ISL 406801 and EPI ISL 412979) were collected from Wuhan.  
312 Additionally, there was another new haplotype (H28, EPI ISL 412980), which  
313 originated from the H14, and was also collected from Wuhan. Although the three  
314 hosts of these two haplotypes didn't directly link with the Huanan seafood market,  
315 one of the hosts (EPI ISL 412979) lived in a residential area about 2 kilometers from  
316 the Huanan seafood market. This also provides the possibility for that host to have  
317 unconsciously indirect contact with the Huanan seafood market. Overall, it cannot  
318 further be proved the SARS-CoV-2 in the Huanan seafood market had been  
319 transmitted from other places or that Huanan seafood market did not host the original  
320 source of SARS-CoV-2. Due to the closing up of the Huanan seafood market on 1<sup>st</sup> of  
321 January, 2020, it made it difficult to determine whether there were intermediate hosts  
322 in the market. To accurately determine whether this market was the central origin,  
323 further collection of early SARS-CoV-2 samples within/around the Huanan seafood  
324 market is needed to provide genomic information and epidemiological survey data  
325 needed in integration analysis. However, based on the facts provided in the current  
326 study undertaken in the early days of the outbreak of the pneumonia infection in  
327 Wuhan, the Huanan seafood market promoted the spread of SARS-CoV-2, and since  
328 then, infected travelers have spread to throughout China and other countries.

329 ***Population size expansion of SARS-CoV-2***

330 As to population demography of SARS-CoV-2, the EBSP results revealed that it  
331 experienced an expansion of effective population size for the last 66 days prior to 11<sup>th</sup>  
332 of February, 2020 (Fig. 5). With the latest one being sampled on 11<sup>th</sup> February 2020,  
333 the first expansion date is estimated to had begun from 7<sup>th</sup> December 2019. The  
334 mismatch in the distribution of the total populations was also clearly shown through  
335 unimodal (Fig. 6) and the neutral tests (Tajima's D and Fu's Fs) revealed statistically  
336 significant population expansion of SARS-CoV-2 (Table S3), which also supported  
337 the EBSP results. In addition, with the change of time, multiple evidence (i.e.  
338 haplotype number, Tajima's D, Fu's FS) indicated a relatively highest population  
339 growth at around 21<sup>st</sup> - 27<sup>th</sup> of January (Table S3). The demographic expansion of  
340 SARS-CoV-2 during this period was consistent with the patients first diagnosed on 8<sup>th</sup>  
341 December, 2019 [4,19] and the previously suggested most recent common ancestor  
342 (TMRCA) dates for SARS-CoV-2 at 6<sup>th</sup> December, 2019) [17]. Additionally,  
343 considering that the virus has longest reported incubation period of up to 24 days [6],  
344 the virus may had first infected humans in mid to late November, which is basically  
345 consistent with the results of previous studies [25]. Therefore, in this study, EBSP  
346 revealed that SARS-CoV-2 experienced an effective population size expansion since  
347 7<sup>th</sup> December 2019, which was also supported by a star-like network, EBSP, the  
348 neutral tests (Fu's and Tajima's D test) and mismatch analysis.

349 **Contributors**

350 Yilin Shu, Tao Pan, and Yongzhong Jiang conceived the research, analyzed the data,  
351 interpreted the results, and wrote the draft manuscript; Bin Fang, Linlin Liu, Xiao Yu  
352 and Xiang Li performed macrogenome sequencing experiments. Xiao Yu and Xiang  
353 Li performed virus isolation experiments. Guojun Ye, Bin Fang, Juan Xu, Ling,

354 Zhang, Yilin Shu and Faxian Zhan collected data. All authors reviewed and approved

355 the final version of the manuscript.

356 **Declaration of interests**

357 We declare no competing interests.

358 **Acknowledgements**

359 We thank scientists and researchers for depositing the whole genomic sequences of

360 Novel Pneumonia Coronavirus (SARS-CoV-2) in the Global Initiative on Sharing All

361 Influenza Data (GISAID) EpiFlu™; and the GISAID database for allowing us access

362 to sequence for non-commercial scientific research. We also are grateful to Dr. Oscar

363 Omondi Donde (Egerton University, Kenya) for polishing the manuscript language.

364 Thanks to Yifa Zhu, Yangyang Tao and Xierong Li (Tianmen Center for Disease

365 Control and Prevention) for collecting samples and thanks to Maoyi Chen, Jie Hu,

366 Chunlin Mao (Jingzhou Center for Disease Control and Prevention) for sampling.

367

368



369 **Reference**

- 370 [1]. Cui J, Li F, Shi ZL. Origin and evolution of pathogenic coronaviruses. Nat  
371 Rev Microbiol. 2019;17(3):181-192.
- 372 [2]. Gorbalenya AE. Severe acute respiratory syndrome-related coronavirus—The  
373 species and its viruses, a statement of the Coronavirus Study Group. BioRxiv.  
374 2020;DOI: <https://doi.org/10.1101/2020.02.07.937862>.
- 375 [3]. Zhu N, Zhang DY, Wang WL, et al. A novel coronavirus from patients with  
376 pneumonia in China, 2019. New Engl J Med. 2020;382(8):727-733.
- 377 [4]. Huang C, Wang Y, Li X, et al. Clinical features of patients infected with 2019  
378 novel coronavirus in Wuhan, China. The Lancet. 2020;395(10223):497-506
- 379 [5]. Lu RJ, Zhao X, Li J, et al. Genomic characterisation and epidemiology of  
380 2019 novel coronavirus: implications for virus origins and receptor binding.  
381 Lancet. 2020; DOI: 10.1016/S0140-6736(20)30251-8.
- 382 [6]. Guan WJ, Ni ZY, Hu Y, et al. Clinical characteristics of 2019 novel  
383 coronavirus infection in China. N Engl J Med. 2020; DOI:  
384 10.1056/NEJMoa2002032.
- 385 [7]. Yang Y, Lu QB, Liu MJ, et al. Epidemiological and clinical features of the  
386 2019 novel coronavirus outbreak in China. MedRxiv.  
387 2020;DOI: <https://doi.org/10.1101/2020.02.10.20021675>.
- 388 [8]. Wassenaar TM, Zou Y. 2019\_nCoV: Rapid classification of betacoronaviruses  
389 and identification of traditional Chinese medicine as potential origin of  
390 zoonotic coronaviruses. Lett Appl Microbiol. 2020; DOI: 10.1111/lam.13285.

- 391 [9]. Zhou P, Yang XL, Wang XG, et al. A pneumonia outbreak associated with a  
392 new coronavirus of probable bat origin. *Nature*. 2020; DOI:  
393 10.1038/s41586-020-2012-7.
- 394 [10]. Wu F, Zhao S, Yu B, et al. A new coronavirus associated with human  
395 respiratory disease in China. *Nature*. 2020; DOI: 10.1038/s41586-020-2008-3.
- 396 [11]. Benvenuto D, Giovannetti M, Ciccozzi A, et al. The 2019-new coronavirus  
397 epidemic: evidence for virus evolution. *J Med Virol*. 2020;  
398 DOI: 10.1101/2020.01.24.915157
- 399 [12]. Ceraolo C, Giorgi FM. Genomic variance of the 2019-nCoV coronavirus. *J*  
400 *Med Virol*. 2020; DOI:10.1002/jmv.25700.
- 401 [13]. Chen LJ, Liu WY, Zhang Q, et al. RNA based mNGS approach identifies a  
402 novel human coronavirus from two individual pneumonia cases in 2019  
403 Wuhan outbreak. *Emerg Microbes Infec*. 2020; 9(1):313-319.
- 404 [14]. Paraskevis D, Kostaki EG, Magiorkinis G, et al. Full-genome evolutionary  
405 analysis of the novel corona virus (2019-nCoV) rejects the hypothesis of  
406 emergence as a result of a recent recombination event. *Infect Genet Evol*. 2020;  
407 79:104212.
- 408 [15]. Chan JF, Kok KH, Zhu Z, et al. Genomic characterization of the 2019 novel  
409 human-pathogenic coronavirus isolated from a patient with atypical  
410 pneumonia after visiting Wuhan. *Emerg Microbes Infec*. 2020;9(1):221-236.
- 411 [16]. Chan JFW, Yuan SF, Kok KH, et al. A familial cluster of pneumonia associated  
412 with the 2019 novel coronavirus indicating person-to-person transmission: a

- 413 study of a family cluster. *Lancet*. 2020; 395(10223):514-523.
- 414 [17]. Li XW, Wang W, Zhao XF, et al. Transmission dynamics and evolutionary  
415 history of 2019-nCoV. *J Med Virol*. 2020; DOI: 10.1002/jmv.25701.
- 416 [18]. Lam TTY, Shum MHH, Zhu HC, et al. Identification of 2019-nCoV related  
417 coronaviruses in Malayan pangolins in southern China. *bioRxiv*. 2020;  
418 DOI: <https://doi.org/10.1101/2020.02.13.945485>.
- 419 [19]. Li Q, Guan XH, Wu P, et al. Early transmission dynamics in Wuhan, China, of  
420 novel coronavirus–infected pneumonia. *New Engl J Med*. 2020;  
421 DOI: 10.1056/NEJMoa2001316.
- 422 [20]. Cohen J. Wuhan seafood market may not be source of novel virus spreading  
423 globally. *Science*. 2020;10.1126/science.abb0611.
- 424 [21]. Marchler-Bauer A, Bo Y, Han LY, et al. CDD/SPARCLE: functional  
425 classification of proteins via subfamily domain architectures. *Nucleic Acids*  
426 *Res*. 2017;45(D1):D200-D203.
- 427 [22]. Katoh K, Standley DM. MAFFT multiple sequence alignment software  
428 version 7: improvements in performance and usability. *Mol Biol Evol*.  
429 2013;30(4):772-780.
- 430 [23]. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and  
431 post-analysis of large phylogenies. *Bioinformatics*. 2014;30:1312–1313.
- 432 [24]. Rambaut A. FigTree v1.4.3; 2016. Available at:  
433 <http://tree.bio.ed.ac.uk/software/figtree/> 2016.
- 434 [25]. Yu WB, Tang GD, Zhang L, et al. Decoding the evolution and transmissions of

- 435 the novel pneumonia coronavirus (SARS-CoV-2) using whole genomic data.  
436 ChinaXiv. 2020;202002.00033v2.
- 437 [26]. Bandelt H, Forster P, Röhl A. Median-joining networks for inferring  
438 intraspecific phylogenies. *Mol Biol Evol.* 1999;16(1):37-48.
- 439 [27]. Librado P, Rozas J. DnaSP v5: a software for comprehensive analysis of DNA  
440 polymorphism data. *Bioinformatics.* 2009;25(11):1451–1452.
- 441 [28]. Tamura K, Stecher G, Peterson D, et al. MEGA6: molecular evolutionary  
442 genetics analysis version 6.0. *Mol Biol Evol.* 2013;30(12):2725–2729.
- 443 [29]. Harpending HC. Signature of ancient population growth in a low-resolution  
444 mitochondrial DNA mismatch distribution. *Hum Biol.* 1994 ;66(4):591-600.
- 445 [30]. Tajima F. The effect of change in population size on DNA polymorphism.  
446 *Genetics.* 1989;123(3):597–601.
- 447 [31]. Fu YX. Statistical tests of neutrality of mutations against population growth,  
448 hitchhiking and background selection. *Genetics.* 1997;147(2):915–925.
- 449 [32]. Drummond AJ, Suchard MA, Xie D, et al. Bayesian phylogenetics with  
450 BEAUti and the BEAST 1.7. *Mol Biol Evol.* 2012;29(8):1969–1973.
- 451 [33]. Liu P, Jiang JZ, Wan XF, et al. Are pangolins the intermediate host of the 2019  
452 novel coronavirus (2019-nCoV) ? bioRxiv. 2020;DOI:  
453 <https://doi.org/10.1101/2020.02.18.954628>.
- 454 [34]. Holshue ML, DeBolt C, Lindquist S, et al. First Case of 2019 Novel  
455 Coronavirus in the United States. *N Engl J Med* 2020;  
456 10.1056/NEJMoa2001191.

**Table 1. Sample information for 10 SARS-Cov-2 infected patients.**

Sample name	Sample type	Collection date	Collection site	Exposure to huanan seafood marke
BetaCoV/Wuhan/HBCDC-HB-02/2019	Bronchoalveolar lavage fluid	30 <sup>th</sup> Dec,2019	Wuhan, Hubei Province	Yes
BetaCoV/Wuhan/HBCDC-HB-01/2019	Bronchoalveolar lavage fluid	30 <sup>th</sup> Dec,2019	Wuhan, Hubei Province	Yes
BetaCoV/Wuhan/HBCDC-HB-03/2019	Bronchial scraping	30 <sup>th</sup> Dec,2019	Wuhan, Hubei Province	Yes
BetaCoV/Wuhan/HBCDC-HB-04/2019	Sputum	30 <sup>th</sup> Dec,2019	Wuhan, Hubei Province	No
BetaCoV/Wuhan/HBCDC-HB-04/2020	Throat swab	18 <sup>th</sup> Jan,2020	Wuhan, Hubei Province	No
BetaCoV/Wuhan/HBCDC-HB-03/2020	Throat swab	18 <sup>th</sup> Jan,2020	Wuhan, Hubei Province	No
BetaCoV/Wuhan/HBCDC-HB-02/2020	Throat swab	17 <sup>th</sup> Jan,2020	Wuhan, Hubei Province	No
BetaCoV/Jingzhou/HBCDC-HB-01/2020	Throat swab and cultured virus	8 <sup>th</sup> Jan,2020	Jingzhou, Hubei Province	Yes
BetaCoV/Wuhan/HBCDC-HB-06/2020	Faeces	7 <sup>th</sup> Feb,2020	Wuhan, Hubei Province	No
BetaCoV/Tianmen/HBCDC-HB-07/2020	Throat swab	8 <sup>th</sup> Feb,2020	Tianmen, Hubei Province	No

459 **Figure Legend**

460 **Fig. 1. Sliding window analyses showing the nucleotide diversity based on**  
461 **alignment of genomes of SARS-CoV-2.** The red line shows the value of nucleotide  
462 diversity ( $\pi$ ) in a sliding window analysis of window size 500 bp with step size 50, the  
463 value is inserted at its mid-point. EPI ISL 402124 sequence was selected as the  
464 reference genome to indicate the mutation position.

465

466 **Fig. 2. Pairwise  $F_{ST}$  among 80 haplotypes of SARS-CoV-2.** Consistent with the  
467 Table S4.

468

469 **Fig. 3. Simplify phylogenetic tree of SARS-CoV-2 form Fig. S1.** The nodal  
470 numbers are ML bootstrap values.

471

472 **Fig. 4. Median-joining network with node sizes proportional to the frequencies of**  
473 **haplotypes in SARS-CoV-2 (A, B and C).** (A) All the individuals collected during  
474 24<sup>th</sup> December to 30<sup>th</sup> December in 2019; (B) All the individuals collected during 24<sup>th</sup>  
475 December in 2019 to 6<sup>th</sup> January in 2020; (C) All the individuals collected before 11<sup>th</sup>  
476 February in 2019. The numbers of mutations separating the haplotypes are shown on  
477 the branches, except for the lower than seven-step mutations. The little red diamond  
478 nodes indicate undetected haplotypes. The sampling areas are indicated by different  
479 colors.

480

481 **Fig. 5. Population fluctuation inferred by Extended bayesian skyline plot (EBSP)**  
482 **of SARS-CoV-2.** The x-axis indicates time in days BP, and the y-axis indicates the  
483 effective population size divided by generation time in units of  $N_{et}$  (the product of  
484 effective population size and generation time in days). The blue areas represent 95 %  
485 highest posterior density. Time is expressed in days.

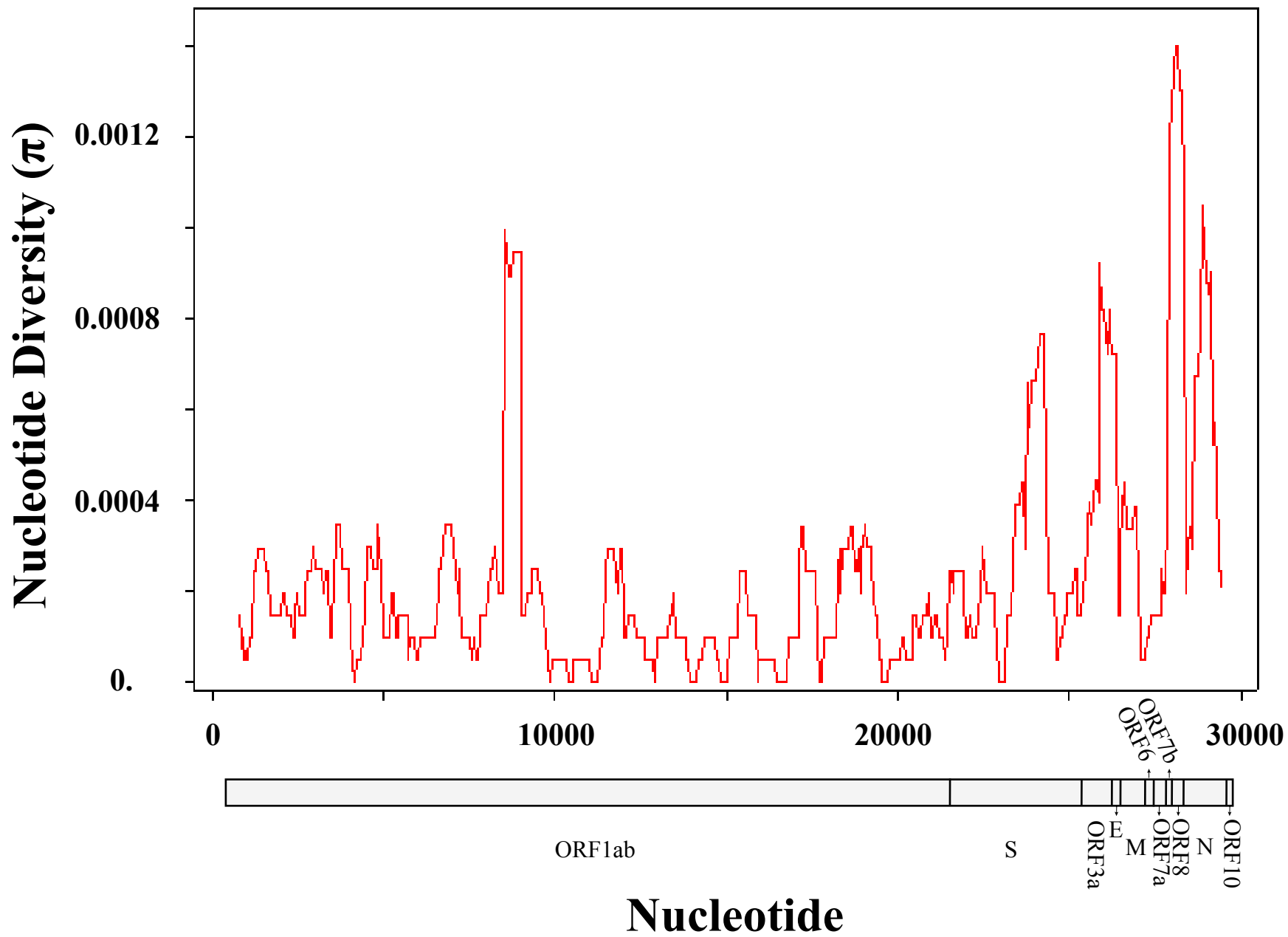
486

487 **Fig. 6: Mismatch distributions analyses for SARS-CoV-2 in different period (A,**  
488 **B, C, D, E, F, G).** (A) All the individuals collected during 24<sup>th</sup> December to 30<sup>th</sup>  
489 December, 2019; (B) 24<sup>th</sup> December in 2019 to 6<sup>th</sup> January in 2020; (C) 24<sup>th</sup>  
490 December in 2019 to 13<sup>th</sup> January in 2020; (D) 24<sup>th</sup> December in 2019 to 20<sup>th</sup>  
491 January in 2020; (E) 24<sup>th</sup> December in 2019 to 27<sup>th</sup> January in 2020; (F) 24<sup>th</sup>  
492 December in 2019 to 3<sup>th</sup> February in 2020; (G) 24<sup>th</sup> December in 2019 to 11<sup>th</sup>  
493 February in 2020. The x coordinate represents the number of differences in each pair  
494 of sequence comparisons; the y coordinate represents the frequencies of pairwise  
495 differences. The blue histogram are the observed frequencies of pairwise divergences  
496 among sequences and the red line refers to the expectation under the model of  
497 population expansion.

498

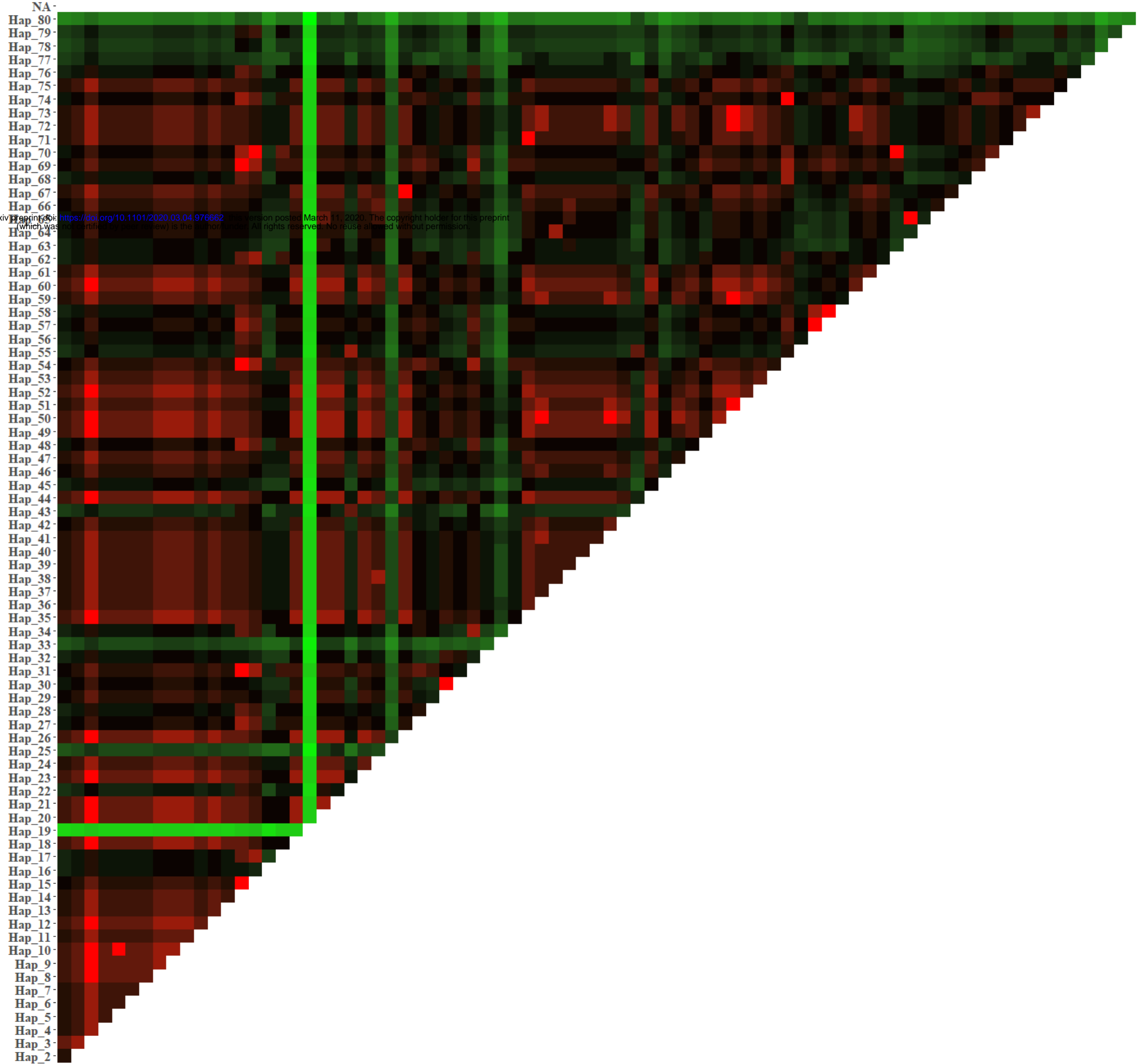
499 **Fig. S1. Phylogenetic tree of SARS-CoV-2.** The nodal numbers are ML bootstrap  
500 values.

501



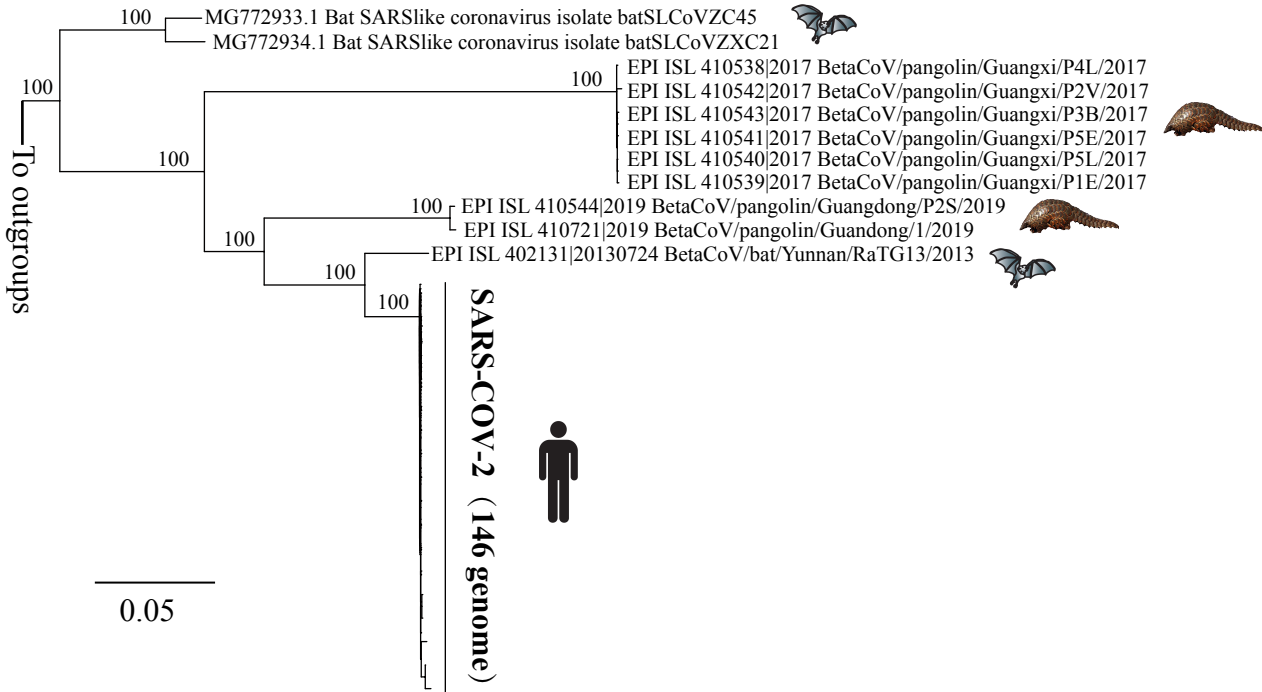


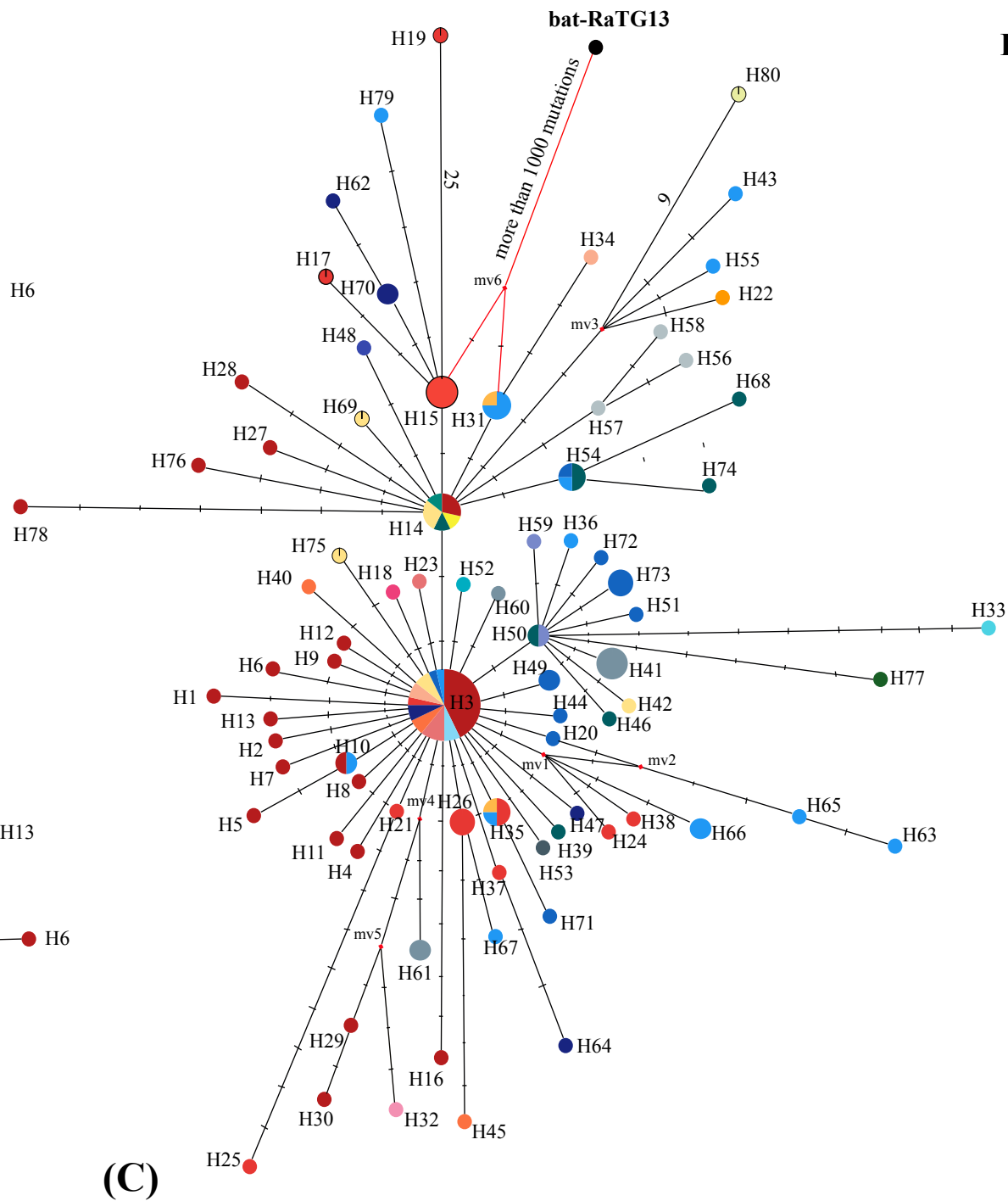
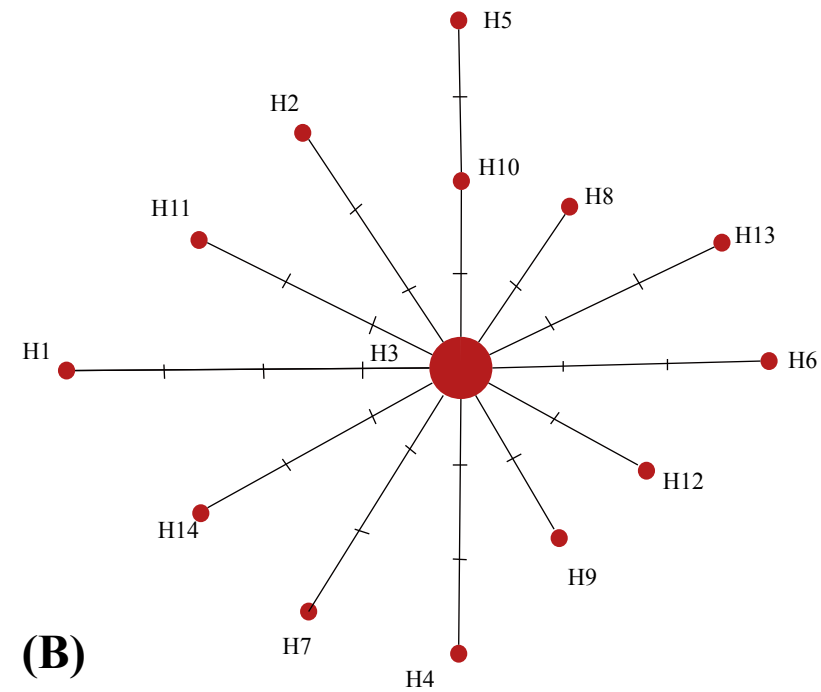
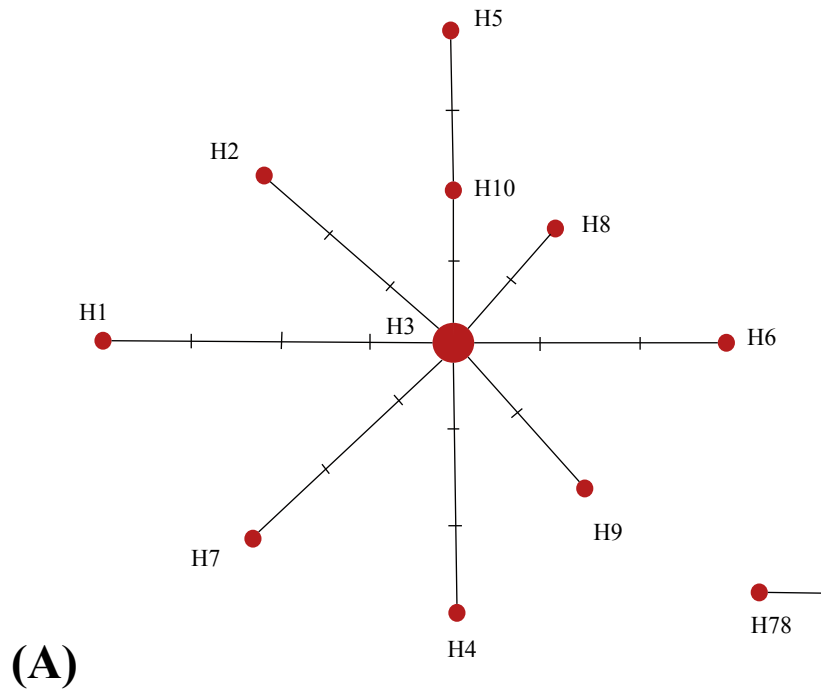
Hap\_1  
Hap\_2  
Hap\_3  
Hap\_4  
Hap\_5  
Hap\_6  
Hap\_7  
Hap\_8  
Hap\_9  
Hap\_10  
Hap\_11  
Hap\_12  
Hap\_13  
Hap\_14  
Hap\_15  
Hap\_16  
Hap\_17  
Hap\_18  
Hap\_19  
Hap\_20  
Hap\_21  
Hap\_22  
Hap\_23  
Hap\_24  
Hap\_25  
Hap\_26  
Hap\_27  
Hap\_28  
Hap\_29  
Hap\_30  
Hap\_31  
Hap\_32  
Hap\_33  
Hap\_34  
Hap\_35  
Hap\_36  
Hap\_37  
Hap\_38  
Hap\_39  
Hap\_40  
Hap\_41  
Hap\_42  
Hap\_43  
Hap\_44  
Hap\_45  
Hap\_46  
Hap\_47  
Hap\_48  
Hap\_49  
Hap\_50  
Hap\_51  
Hap\_52  
Hap\_53  
Hap\_54  
Hap\_55  
Hap\_56  
Hap\_57  
Hap\_58  
Hap\_59  
Hap\_60  
Hap\_61  
Hap\_62  
Hap\_63  
Hap\_64  
Hap\_65  
Hap\_66  
Hap\_67  
Hap\_68  
Hap\_69  
Hap\_70  
Hap\_71  
Hap\_72  
Hap\_73  
Hap\_74  
Hap\_75  
Hap\_76  
Hap\_77  
Hap\_78  
Hap\_79



bioRxiv preprint doi: <https://doi.org/10.1101/2020.03.04.976662>; this version posted March 11, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

$P_i$   
 $10^{-3}$   
 $10^{-3.5}$   
 $10^{-4}$





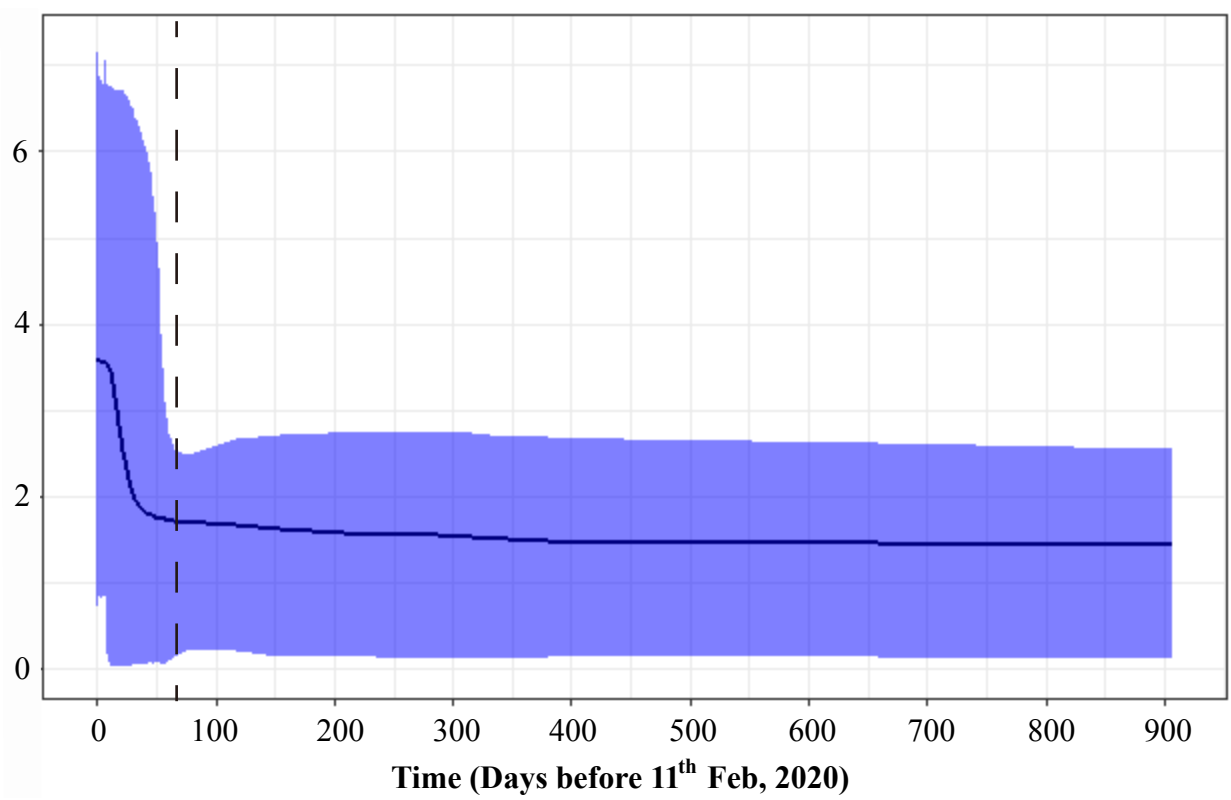
## Legends

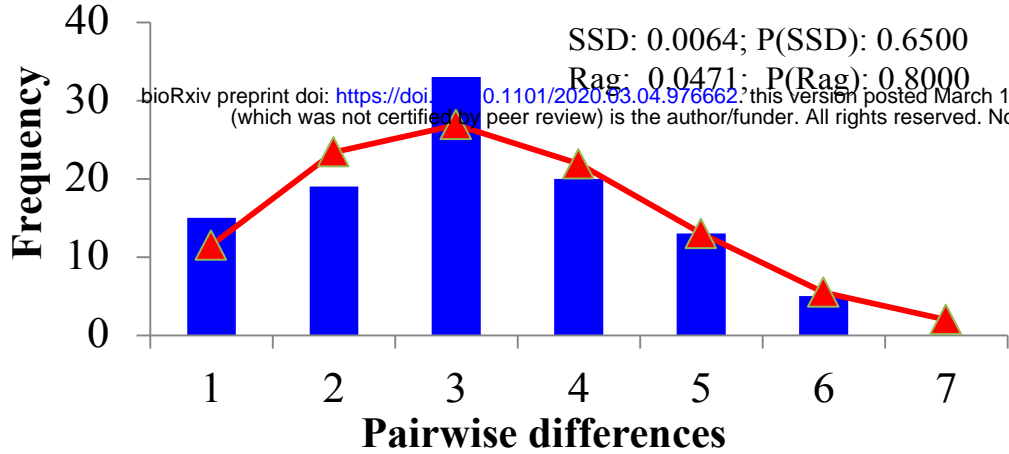
### China

- Hubei
- Guangdong
- Zhejiang
- Jiangxi
- Shandong
- Jiangsu
- Chongqing
- Sichuan
- Fujian
- Yunnan
- Taiwan
- Beijing

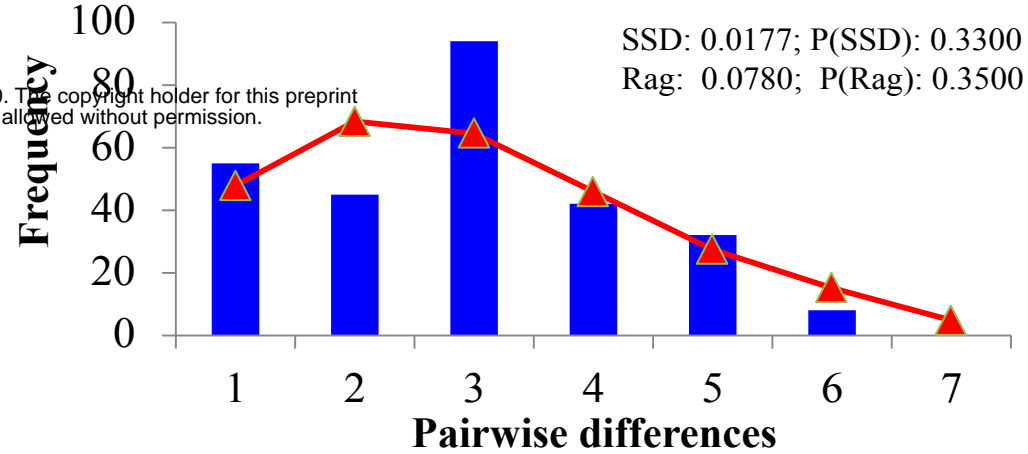
### Other Countries

- Japan
- Korea
- Italy
- Singapore
- United States
- Nonthaburi
- Australia
- Cambodia
- South Korea
- Sweden
- Belgium
- Nepal
- Germany
- France
- England

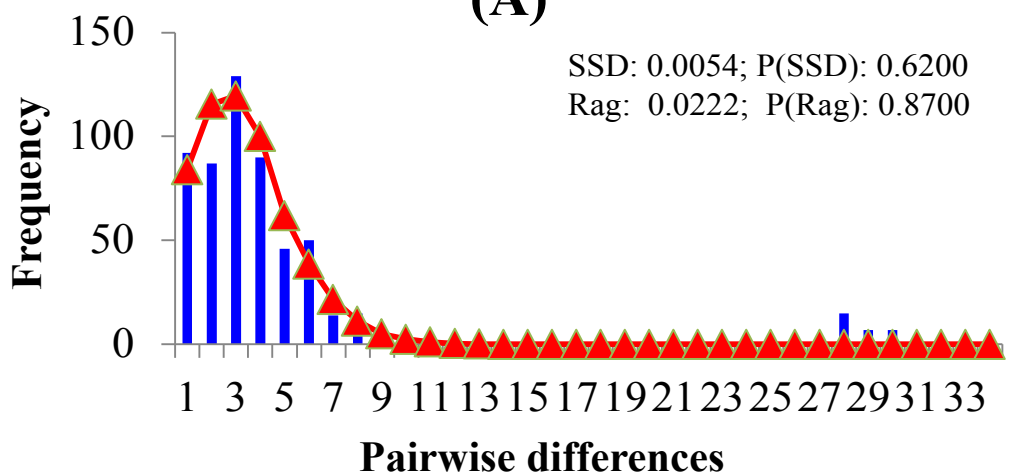




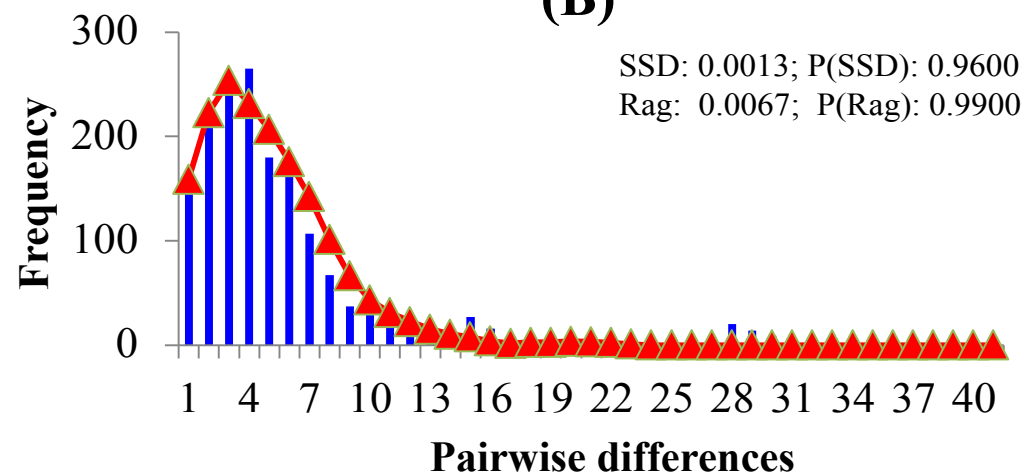
(A)



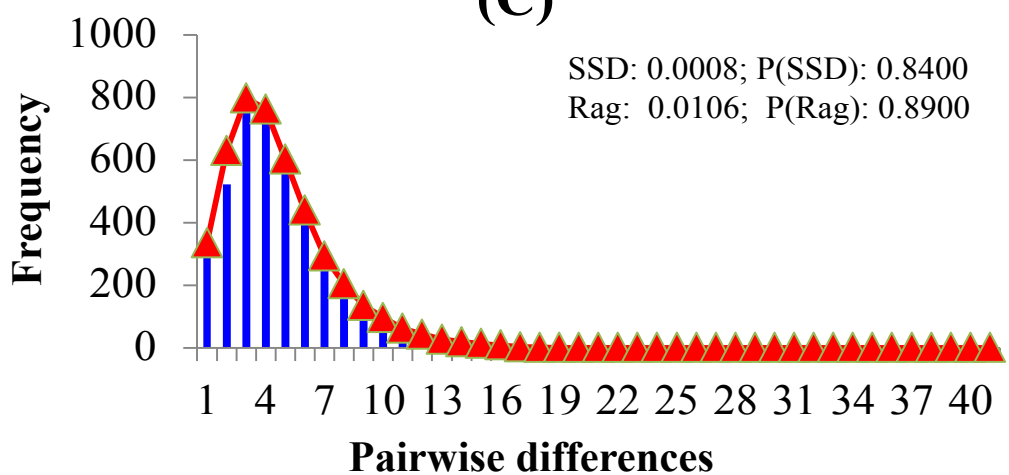
(B)



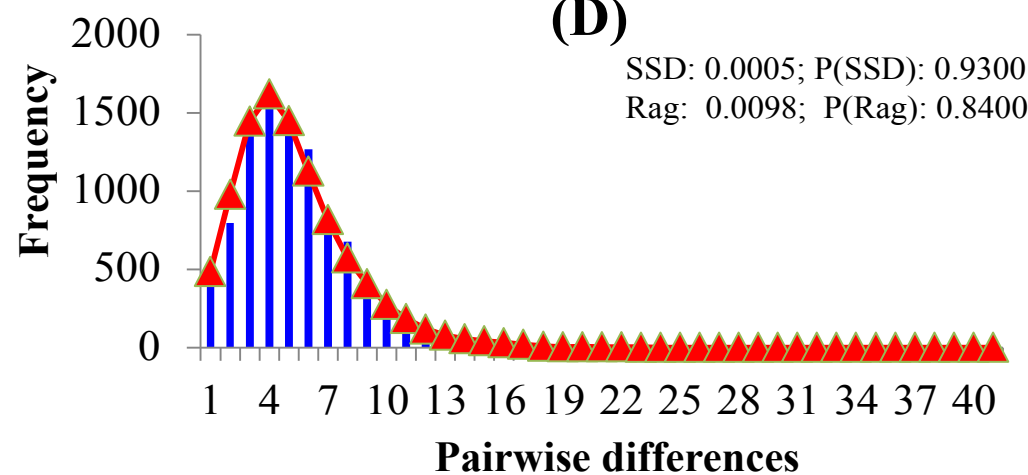
(C)



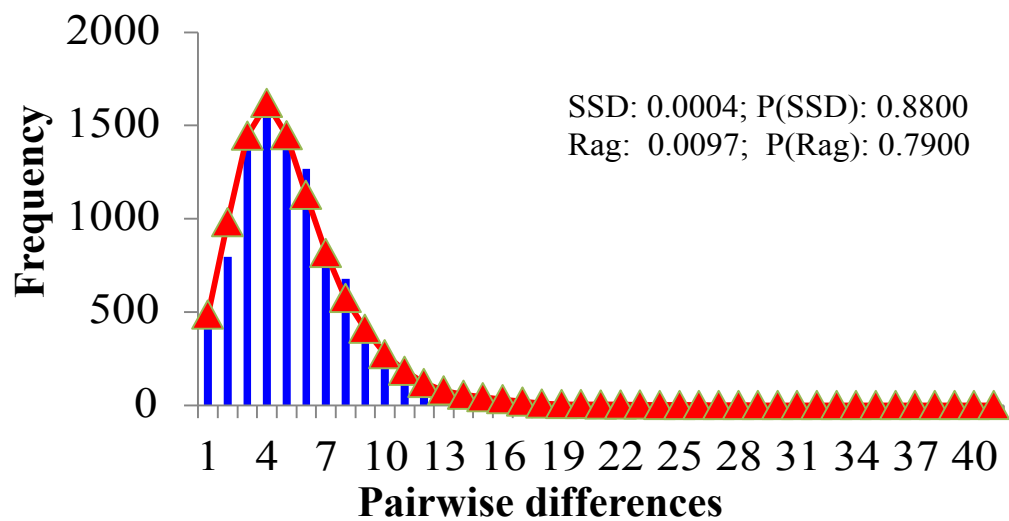
(D)



(E)



(F)



(G)