

1                                   **Identification of super-transmitters of SARS-CoV-2**

2

3                                   Xuemei Yang<sup>1#</sup>, Ning Dong<sup>1#</sup>, Edward Wai-Chi Chan<sup>2</sup>, Sheng Chen <sup>1\*</sup>

4

5       <sup>1</sup>Department of Infectious Diseases and Public Health, Jockey Club College of Veterinary  
6       Medicine and Life Sciences, City University of Hong Kong, Kowloon, Hong Kong

7

8       <sup>2</sup>State Key Lab of Chemical Biology and Drug Discovery, Department of Applied Biology and  
9       Chemical Technology, The Hong Kong Polytechnic University, Hung Hom, Kowloon, Hong  
10      Kong;

11

12      #contribute equally to the work.

13

14      \*Corresponding author: Sheng Chen, City University of Hong Kong, Kowloon, Hong Kong;

15      Email: [shechen@cityu.edu.hk](mailto:shechen@cityu.edu.hk)

16

17

18      **Keywords:** SARS-CoV-2, COVID-19, Origin, phylodynamic analysis, Super-transmitter

19

20 **Abstract**

21 A newly emerged coronavirus, SARS-CoV-2, caused severe outbreaks of pneumonia in  
22 China in December 2019 and has since spread to various countries around the world. To  
23 probe the origin and transmission dynamics of this virus, we performed phylodynamic  
24 analysis of 247 high quality genomic sequences of viruses available in the GISAID platform  
25 as of March 05, 2020. A substantial number of earliest sequences reported in Wuhan in  
26 December 2019, including those of viruses recovered from the Huanan Seafood Market  
27 (HNSM), the site of the initial outbreak, were genetically diverse, suggesting that viruses of  
28 multiple sources were involved in the original outbreak. The viruses were subsequently  
29 disseminated to different parts of China and other countries, with diverse mutational profiles  
30 being recorded in strains recovered subsequently. Interestingly, four genetic clusters defined  
31 as Super-transmitters (STs) were found to become dominant and were responsible for the  
32 major outbreaks in various countries. Among the four clusters, ST1 is widely disseminated in  
33 Asia and the US and mainly responsible for outbreaks in the states of Washington and  
34 California in the US as well as those in South Korea at the end of February and early March,  
35 whereas ST4 contributed to the pandemic in Europe. Each ST cluster carried a signature  
36 mutation profile which allowed us to trace the origin and transmission patterns of specific  
37 viruses in different parts of the world. Using the signature mutations as markers of STs, we  
38 further analysed 1539 genome sequences reported after February 29, 2020. We found that  
39 around 90% of these genomes belonged to STs with ST4 being the dominant one and their  
40 contribution to pandemic in different continents were also depicted. The identification of  
41 these super-transmitters provides insight into the control of further transmission of SARS-  
42 CoV-2.

## 43 **Introduction**

44 A number of newly emerged coronaviruses such as the highly pathogenic severe acute  
45 respiratory syndrome coronavirus (SARS-CoV) and Middle East respiratory syndrome  
46 coronavirus (MERS-CoV) have caused serious respiratory and intestinal infections in human  
47 within the past two decades [1]. In December 2019, another new coronavirus, SARS-CoV-2,  
48 has emerged and caused outbreaks of lower respiratory tract infections, often with poor clinical  
49 outcome, in Wuhan, China. The virus, which has since spread to other cities in China and  
50 various countries worldwide [2], exhibited a high potential to undergo human-to-human  
51 transmission [3]. As of March 26, 2020, a total of 692 thousand infections were documented  
52 worldwide, among which 68 thousand occurred in China ([https://www.gisaid.org/epiflu-](https://www.gisaid.org/epiflu-applications/global-cases-betacov/)  
53 [applications/global-cases-betacov/](https://www.gisaid.org/epiflu-applications/global-cases-betacov/)). As a result, WHO declared the risk of SARS-CoV-2 as  
54 “Very High” in China ([https://www.who.int/docs/default-source/coronaviruse/situation-](https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200202-sitrep-13-ncov-v3.pdf)  
55 [reports/20200202-sitrep-13-ncov-v3.pdf](https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200202-sitrep-13-ncov-v3.pdf)).

56  
57 The genomic characteristics of SARS-CoV-2 have been elucidated using phylogenetic,  
58 structural and mutational analyses by scientists across the globe [4]. High-throughput  
59 sequencing revealed that SARS-CoV-2 was a novel betacoronavirus which resembled SARS-  
60 CoV at around 79.5% sequence identity [5, 6]. A recent study indicated that SARS-CoV-2 was  
61 96% identical to a bat coronavirus RaTG13 (accession: MN996532) at the genomic level,  
62 suggesting that bat might be a natural host of SARS-CoV-2 [7]. GISAID is a platform for  
63 sharing genetic data of influenza. Currently, a rapidly increasing number of SARS-CoV-2  
64 genomic sequences are being deposited into this database from laboratories around the world  
65 [8]. On the other hand, some recent studies also suspected that Malayan pangolins (*Manis*  
66 *javanica*) could be the intermediate host of this new coronavirus, since the amino acid sequence  
67 of the S protein of coronaviruses derived from Malayan pangolins illegally imported to

68 Guangdong Province of China, as well as coronaviruses harboured by pangolins in Guangxi  
69 province of China, exhibited very high homology with the S protein sequence of SARS-CoV-  
70 2, even though the overall homology between SARS-CoV-2 and RaTG13 is still the highest[9].  
71 However, due to the inability to detect or isolate SARS-CoV-2 from pangolins in the Wuhan  
72 Huanan seafood wholesale market, the site in which the first batch infected patients had  
73 commonly visited, the theory of pangolins being the culprit of the Wuhan pneumonia outbreak  
74 is not substantiated. The intermediate host for SARS-CoV-2 therefore remains a mystery. In  
75 fact, it remains unclear if the Huanan Seafood Wholesale Market is the origin of this outbreak  
76 as some of the earliest cases were confirmed to have no linkage with this market. It is urgent  
77 to identify the source(s) of viruse(s) that caused this outbreak to design more effective control  
78 measures to stop the continuous worldwide transmission of these highly contagious viruses.  
79 With more sequences released, it is very important to provide more insights of this virus  
80 through indepth sequence analysis. One recent study has analysed over 100 available genome  
81 sequences and revealed that sequences belonging to different genetic clusters have evolved  
82 [10]. In this study, we retrieved and analysed the publicly shared genome sequences as of  
83 March 25, 2020 to investigate the genetic diversity and phylodynamics of these SARS-CoV-2  
84 viruses. We identified notable four clusters of genomes with high transmission and mutation  
85 rate and have become the most dominant viruses in the later stage of transmission. Results in  
86 this study should provide insight into the control of further transmission of SARS-CoV-2.

87

88

## 89 **Materials and Methods**

### 90 **Sequence analysis, alignment and mutation identification**

91 A total of 343 full-length SARS-CoV-2 genomes available in the GISAID platform  
92 (<https://platform.gisaid.org/>) as of March 5, 2020 were downloaded [8]. A total of 247  
93 sequences with high sequence quality as noted in the GISAID database were included for  
94 further analysis after removing sequences containing little temporal signal and thus are not  
95 unsuitable for inference using phylogenetic molecular clock models. Information regarding the  
96 date and country of isolation were also retrieved from the GISAID platform. The annotated  
97 reference genome sequence of the SARS-CoV-2 isolate Wuhan-Hu-1 (accession:  
98 NC\_045512.2) was downloaded from the NCBI GenBank database. All genomes were  
99 annotated by GATU Genome Annotator [11] using the SARS-CoV2 isolate Wuhan-Hu-1  
100 (NC\_045512.2) as reference [12]. Nucleotide and amino acid mutations of all genome and  
101 separate proteins were analyzed by blast (<https://blast.ncbi.nlm.nih.gov/>) using the sequence of  
102 strain Wuhan-Hu-1 as reference.

103

### 104 **Phylogenetic analysis**

105 Global genomic surveillance of SARS-CoV-2 was implemented by means of an automated  
106 phylogenetic analysis pipeline using Nextstrain, which generates an interactive visualization  
107 integrating a phylogeny with sample metadata such as geographic location or host age [13].  
108 The pipeline involved the sequence alignment module with MAFFT [14], phylogenetic  
109 analysis with IQ-TREE [15], maximum-likelihood phylodynamic analysis with Treetime [16],  
110 identification of nucleotide and amino acid mutations with Augur, and result visualization with  
111 Auspice [13]. The outputs were edited by Inkscape 0.91 [17].

112

### 113 **Phylogenetic analysis**

114 Alignment of the complete genome sequences was conducted with MAFFT v7.310 [14].  
115 Phylogenetic tree of all SARS-CoV-2 was built with RAxML version 8.2.4 [18]. The tree was  
116 edited by iTOL [19].

117

### 118 **Quick identification of the types of SARS-CoV-2 genomes in the database**

119 All complete genomes available as March 28 on the GISAID database were downloaded.  
120 Single Nucleotide Polymorphisms (SNPs) calling were performed by Snippy  
121 (<https://github.com/tseemann/snippy>) using Wuhan-Hu-1 as reference. Super-transmitter  
122 clusters were classified according relative variants. A total of 1956 qualified genomes  
123 submitted after February 29, 2020 were included.

124

## 125 **Results**

### 126 **Phylogenetics analysis of genome sequences of SARS-CoV-2 strains collected worldwide**

127 To trace the evolutionary process and identify the common ancestor of 247 strains of SARS-  
128 CoV-2 collected worldwide, root-to-tip regression scatter plots was conducted among all  
129 SARS-CoV-2 genomes, with  $R^2$  being found to be 0.23, suggesting that these 247 viral  
130 sequences shared a common recent ancestor (**Fig 1a**). The date of the most recent common  
131 ancestor (tMRCA) of all reported SARS-CoV-2 viruses was 2019-Nov-12, suggesting that this  
132 virus emerged recently (**Fig 1a**). A total of 379 nucleotide mutations were identified among  
133 these 247 sequences based on the sequence alignment, among which G<sup>11083</sup>T (n=5), T3G (n=3),  
134 G<sup>29864</sup>A (n=3), C<sup>29870</sup>A (n=3), A1T (n=2), A4T (n=2), T<sup>4402</sup>C (n=2), G<sup>5062</sup>T (n=2), T<sup>18603</sup>C (n=2)  
135 and G<sup>22661</sup>T (n=2) were the most homoplasic mutations (Fig 2, supplementary file 1). A total  
136 of 147 strains were found to contain single amino acid change, with the majority of such  
137 changes being located within ORF1ab (n=104). The L<sup>3606</sup>F change was detected in two viral  
138 sequences, while other mutations occurred only once. Mutations that result in amino acid  
139 changes include single substitution in the S protein (n=19, D<sup>614</sup>G, L<sup>752</sup>F, F<sup>32</sup>I, H<sup>655</sup>Y, V<sup>483</sup>A,  
140 F<sup>157</sup>L, V<sup>615</sup>L, K<sup>202</sup>N, S<sup>939</sup>F, F<sup>797</sup>C, A<sup>93</sup>0V, R<sup>408</sup>I, V<sup>367</sup>F, Q<sup>409</sup>E, S<sup>254</sup>F, A435S, D<sup>1146</sup>E, S<sup>247</sup>R and  
141 P<sup>1143</sup>L), ORF3a (n=8, E<sup>191</sup>G, G<sup>76</sup>S, K<sup>61</sup>N, V<sup>259</sup>L, T<sup>176</sup>I, L<sup>140</sup>V, T<sup>269</sup>M and V<sup>88</sup>L), N protein (n=6,  
142 K<sup>247</sup>I, S<sup>194</sup>L, P<sup>46</sup>S, S<sup>327</sup>L, E<sup>378</sup>Q and D<sup>343</sup>V), ORF8 (n=4, T<sup>11</sup>I, L<sup>84</sup>S, S<sup>97</sup>N and S<sup>67</sup>F), ORF7a  
143 (n=3, P<sup>34</sup>S, Q62\* and H<sup>73</sup>Q), ORF10 (n=2, P<sup>10</sup>S and I<sup>13</sup>M) and E protein (n=1, S<sup>6</sup>L) (**Fig 1b**).  
144 Identification of single amino acid substitutions in SARS-CoV-2 isolates consistently showed  
145 that these isolates shared a recent common ancestor but entered diverse evolution paths. The  
146 estimated substitution rate of SARS-CoV-2 was 8.90e-04 subs/site/year, which was similar to  
147 that of other RNA viruses including SARS-CoV, Ebola virus, Zika virus, and others, which  
148 was found to be at ~ 1e-3 subs/site/year ([http://virological.org/t/phylogenetic-analysis-93-  
149 genomes-15-feb-2020/356](http://virological.org/t/phylogenetic-analysis-93-genomes-15-feb-2020/356)). Based on this mutation rate, a genome of 29kb of SARS-CoV-2 will

150 end up with ~26 mutations per genome per year, suggesting that within the two months' study  
151 period, the number of mutations in each genome should not exceed five if all test isolates  
152 emerged as a result of natural evolution of a single SARS-CoV-2 strain.

153

#### 154 **Multiple origins of SARS-CoV-2 in Wuhan, China**

155 To shed light on the evolution trend of SARS-CoV-2, we analysed the time-dependent changes  
156 in mutation profiles of the test strains in detail. A total of 16 viral genomes collected before  
157 Jan. 01, 2020, were included (**Table 1**). All of these 16 genomes were obtained from Wuhan,  
158 with half of them from Huanan Seafood Wholesale Market (HNSM) where the original  
159 outbreak occurred. Six genomes contained identical sequences, four of which belonged to  
160 isolates obtained from HNSM. Compared to these six viral genomes, others displayed various  
161 mutation profiles which comprised 1 to 6 mutations in the genomes. We therefore set these six  
162 genomes as reference genome for subsequent analyses. Two earliest viral genomes reported on  
163 Dec. 24 and 26 were found to harbour two and three mutations when compared to the reference  
164 viral genome, respectively. Four viral genomes from HNSM also contained two mutations with  
165 different profiles, suggesting that the original SARS-CoV-2 strain might have been circulating  
166 in HNSM for a certain period of time and underwent mutational changes in different  
167 intermediate hosts (**Table S1**). These observations suggested that HNSM was not the only  
168 origin of the COVID-19 outbreak, instead the market might only serve as a medium in which  
169 transmission of this virus to human first occurred. The original virus seemed to have  
170 transmitted to various provinces in China subsequently, including Guangdong, Zhejiang,  
171 Anhui, Jiangshu and Chongqing, and then to other countries including Japan, Taiwan, Thailand  
172 and USA in the following month (January 2020). The viral genome reported initially from USA  
173 were those of viruses recovered from the patients in Princess Diamond Cruise, confirming that  
174 the original virus was the one that caused the outbreak in this cruise; such view is consistent

175 with the finding that identical genomes were reported in Japan, where the cruise ship was  
176 docked. A total of 26 out of the 247 genome sequences tested contain one mutation. Unless  
177 isolated from the same location, most of these genomes exhibit unique mutational profile. Five  
178 sequences from the Princess Diamond cruise ship were found to exhibit unique mutation  
179 profiles, thus further suggesting that random mutations occurred during viral evolution. It  
180 should be noted that these genome sequences were also reported in Wuhan, other parts of China  
181 and various other countries, confirming that the transmission of the original virus to different  
182 parts of the world was accompanied by active but random mutational changes during the  
183 process (**Table S1**).

184

#### 185 **Phylogenetic analysis of genome sequences of SARS-CoV-2**

186 Phylogenetic analysis of the 247 SARS-CoV-2 genomes was also performed, with results  
187 showing that such viral genomes exhibited highly diverse genetic profiles and that random  
188 mutations occurred during the evolutionary process within the first two months. Interestingly,  
189 four clusters of genome sequences were observed among the 247 genomes, with the rest  
190 exhibiting more diverse profiles. These results were consistent with the data of maximum-  
191 likelihood phylodynamic analysis shown in Figure 1. Comparison of the mutation profile of  
192 each cluster enabled us to discover that all viral genomes in the same cluster were derived from  
193 one parental viral genome sequence which bears a signature mutation profile, as such profile  
194 could be identified in all offsprings (**Fig 2**). The first cluster contained two mutations, C<sup>8782</sup>T  
195 and T<sup>28144</sup>C; the second cluster contained the mutation G<sup>26144</sup>T; the third cluster contained the  
196 mutation G<sup>11083</sup>T; the fourth cluster contained three mutations, C<sup>241</sup>T, C<sup>3037</sup>T and A<sup>23403</sup>G.  
197 Tracing the changes in mutation profiles of these viral genomes over time allowed us to  
198 visualize the transmission and evolution dynamics of SARS-CoV-2. Since viruses of all of

199 these four clusters exhibited very high potential to undergo global transmission, we define  
200 viruses in these four clusters as super-transmitter cluster 1 (ST1), 2 (ST2), 3(ST3) and 4(ST4).

201

### 202 **Evolution and transmission of super-transmitter cluster 1 (ST1)**

203 ST1 carried the signature mutation profile of C<sup>8782</sup>T and T<sup>28144</sup>C. The C<sup>8782</sup>T change is a silent  
204 mutation, whereas T<sup>28144</sup>C is associated with the amino acid substitution L<sup>84</sup>S in the Orf8ab  
205 protein. The ST1 viruses were transmitted very efficiently and a total of 85 out of 247 (34%)  
206 genome sequences belonging to this cluster as of March 03, 2010. The earliest sequence of in  
207 this cluster was reported in Wuhan, China on Jan 05, 2020 and seven were subsequently  
208 reported in January and February in different parts of China and Australia, suggesting that  
209 widespread transmission of this cluster of viruses occurred (**Table 2**). The viruses in ST1 were  
210 mainly transmitted among Asian countries area such as China, Japan, South Korea, Taiwan  
211 and Singapore, as well as North America, in particular the states of California and Washington  
212 in USA (**Table 2**). The viruses in ST1 were also found to be able to rapidly mutate along the  
213 transmission paths. Three genome sequences that were reported in Australia, Vietnam and USA  
214 on Feb. 28, 24 and Mar. 03, 2020 respectively, were found to harbour a total of 11 mutations.  
215 An additional nine mutations were acquired by the parental virus within 50 days (from Jan 05  
216 to Feb. 24, 2020), with a mutation rate of 2.3e-3 subs/site/year (29kb genome size), which was  
217 much higher than the predicted mutation rate of SARS-CoV-2 (4.057 e-4 subs/site/year) and  
218 other coronaviruses such as SARS-CoV-1 and MERS virus. Among viral genomes in this  
219 cluster, 43 out of 85 genomes exhibited five or more mutations (**Table 2**).

220

221 Detailed analysis of mutation profiles of the genome sequences in ST1 enables us to trace the  
222 evolution routes of these viruses in specific region. In Washington State, USA, a genome  
223 sequence with three mutations, C<sup>18060</sup>T, C<sup>8782</sup>T and T<sup>28144</sup>C, was reported on Jan 25, 2020. A

224 virus carrying these three mutations was reported in Fujian, China on Jan 21, 2020, suggesting  
225 that the virus might have originated from Fujian, China (**Table 2**). This virus was then further  
226 transmitted in the Washington area and continued to acquire mutations. Twelve genome  
227 sequences reported between March 01 - 05, 2020 in the states of Washington and California  
228 contained 6 to 11 mutations. The data provided direct evidence of active evolution that results  
229 in a large number of mutational changes during the process of transmission of a single virus  
230 within a short period of time. In addition, a genome sequence which has two additional  
231 mutations when compared to the original virus in this cluster, but were different from those in  
232 genome sequences in Washington, was reported in Sichuan, China, suggesting that the same  
233 parental virus was also transmitted across China during this period (**Table 2**).

234

### 235 **Evolution and transmission of super-transmitter 2 (ST2)**

236 ST2 carried the signature mutation G<sup>26144</sup>T, which resulted in the G<sup>251</sup>V amino acid substitution  
237 in Orf 3 protein of SARS-CoV-2. The first viral genome in this cluster was reported on Jan. 25,  
238 2020 in Australia and a total of 28 out of the 247 (11.2%) sequences tested were found to  
239 belong to ST2. The parental virus had acquired different mutations and had been disseminated  
240 to various Asian countries, North America (USA), Europe, South America (Brazil) and  
241 transmission in Australia. Viruses in this cluster seemed to be extensively transmitted by the  
242 end of January and lasted till early February. By the end of February, however, transmission  
243 efficiency of such viruses seemed to have dropped, as only 4 out of 28 sequences reported during  
244 the period Feb. 26 to Mar. 03, 2020 belong to this cluster. Viruses in this cluster were also  
245 found to have significantly mutated. Examples are three genomes reported from Switzerland  
246 on Feb. 29, 2020, in which 9, 10 and 11 mutations were identified respectively (**Table 3**). Our  
247 data showed that as many as eight additional mutations were acquired by the parental virus  
248 within 30 days (from Jan 28 to Feb. 29, 2020), representing a mutation rate of 3.3e-3

249 subs/site/year (29kb genome size), which was much higher than the predicted mutation rate of  
250 SARS-CoV-2 ( $8.0e-4$  subs/site/year).

251

### 252 **Evolution and transmission of super-transmitter 3 (ST3)**

253 ST3 carried the signature mutation G<sup>11083</sup>T, which caused the L<sup>3606</sup>F amino acid substitution in  
254 the Orf 1 protein of SARS-CoV-2. The first viral genome in this cluster was reported on Jan.  
255 18, 2020 in Chongqing, China and a total of 22 (9%) genome sequences were reported so far.  
256 It has since been transmitted to some Asian countries including Singapore, as well as Japan,  
257 Europe, USA and Australia (**Table 4**). Like ST1 and ST2, viruses in this cluster were also  
258 found to mutate efficiently, with one genome reported on Feb. 27, 2020 from Washington,  
259 USA, carrying 12 mutations. Our data showed that a total of eleven mutations were acquired  
260 by the parental virus within a 40 days period (from Jan 18 to Feb. 27, 2020), with a mutation  
261 rate of  $2.8e-3$  subs/site/year (29kb genome size), which was again much higher than the  
262 predicted mutation rate of SARS-CoV-2 ( $8.0e-4$  subs/site/year). Curiously, there is no virus of  
263 this cluster being reported in Iran, a country with one of the highest incidence of SARS-CoV-  
264 2 infections. However, two genome sequences from Australia, which belong to viruses  
265 recovered from patients with travel history to Iran, were reported, suggesting that this cluster  
266 of virus might also contribute to the outbreaks in Iran. In addition, the first genome sequence  
267 from Brazil, which might have originated from Italy, also belonged to this cluster (**Table 4**).

268

### 269 **Evolution and transmission of super-transmitter 4 (ST4)**

270 ST4 carried a signature mutation profile that consists of three mutations: C<sup>241</sup>T, C<sup>3037</sup>T and  
271 A<sup>23403</sup>G. The C<sup>241</sup>T and C<sup>3037</sup>T changes are silent mutations, whereas A<sup>23403</sup>G results in the  
272 D<sup>614</sup>G substitution in the spike (S) protein of SARS-CoV-2. ST4 viruses were found to be  
273 transmitted only in Europe, with the exception of one genome from Mexico with travel history

274 from Italy, and contributed to the current explosive increase in incidence of COVID-19 in  
275 Europe (**Table 5**). Genomes in ST4 were reported more recently, from the end of February to  
276 early March. A total of 21 out of 247 (8.4%) genome sequences have been reported so far.  
277 Among the 247 genome sequences of the four clusters of super-transmitters, there is no genome  
278 containing either one of these three mutations or a combination of two of these three mutations,  
279 suggesting that the parental viral genome of ST4 could not be identified. The first virus of this  
280 cluster was reported on Jan 28, 2020 in Germany. This virus acquired another silent mutation,  
281 C<sup>14408</sup>T, and further spread to other countries in Europe. This virus was also found to mutate  
282 efficiently, with three genomes reported on Feb. 29, 2020 from Switzerland carrying 9, 10 and  
283 11 mutations respectively. Eight additional mutations were acquired by the parental virus  
284 within 30 days (from Jan 28 to Feb. 29, 2020), with a mutation rate of 3.3e-3 subs/site/year  
285 (29kb genome size), which was much higher than the predicted mutation rate of SARS-CoV-2  
286 (8.0e-4 subs/site/year). At a later stage, ST4 viruses became the most efficient in causing  
287 transmission in Europe. Among the 28 genomes reported between Feb. 20 to March 03, 2020,  
288 20 (71%) belonged to this transmitter (**Table 5**).

289

## 290 **Temporal and spatial distribution of super-transmitters of SARS-CoV-2**

291 To better understand the temporal and spatial distribution of these super-transmitters, we plot  
292 variation in the types of genome sequences recovered from different continents against time.  
293 The original viruses were transmitted in the first week before the emergence of these super-  
294 transmitters. ST1 was the first batch of viruses that emerged and dissemination continued  
295 throughout the study period. Other STs emerged at different time points and transmission also  
296 peaked at different dates. Transmission of ST2 and ST3 mainly occurred between mid January  
297 to mid February. Transmission of ST4 viruses mainly began at the end of February. Viruses of  
298 the four clusters exhibited much higher mutation rate than those which exhibited diverse

299 genetic profiles and could not be allocated into specific genetic cluster when compared to the  
300 original genome (**Fig 3a**). ST1 viruses were those which were disseminated extensively in  
301 China, in particular in the later stage of the outbreak (**Fig 3b**). ST1, ST2 and ST3 were prevalent  
302 in other Asian countries (**Fig 3c**). All the four clusters were involved in the outbreaks in Europe  
303 in the early stage, but ST4 was the cluster that eventually transformed the outbreaks in Europe  
304 to the pandemic level (**Fig 3d**). In Oceania, ST1 was involved mainly in the early stage of  
305 outbreak, yet ST2 became dominant at the later stage (**Fig 3e**). ST1 and other types of viruses  
306 were the major transmitters in the US. ST1 was shown to be transmitted mainly in the states of  
307 Washington and California, whereas the other types were mainly transmitted in other states,  
308 (**Fig 3f, Table S1**).

309

### 310 **Distribution of different types of most recent SARS-CoV-2 in different parts of the** 311 **world**

312 Upon finishing our manuscript, we went to check the available genome sequences in the  
313 database and found a rapid increase of numbers of sequences. A total of 1539 genome  
314 sequences reported after February 29, 2020 were included for a quick analysis to identify the  
315 type of these most recent genomes. As shown in Table 6, most of the genomes were reported  
316 from USA (968 / 63%) and Europe (441 / 29%), where the pandemics were the most server.  
317 It is good to see some genomes from Africa (20 / 1%) and South America (23 / 2%), which  
318 were minimally reported before March 01, 2020. Among these genomes, 89% of the genomes  
319 belonged to ST1-4 with ST4 being the most dominant (56%), while the original derivatives  
320 accounted for only 11%, which were mainly reported in UK and Netherland. In Africa, ST4  
321 (18/20, 90%) was the major type with some of the cases showing travel history to Europe; in  
322 Asia, the major types became ST3 (17/33, 52%) and ST4 (16/33, 48%); in Europe, all types  
323 were presence with ST4 being the dominant one (668/968, 69%); all the types were reported  
324 in the US with ST1 (282/441, 62%) and ST4 (137/441, 14%) being the dominant; in Canada,

325 all types except for ST1 were present; in Oceania, all four STs were present with ST3 being  
326 the dominant; in South America, ST1 and ST4 were the most dominant types (**Table 6**).

## 327 **Discussion**

328 We conducted detailed and comprehensive analyses of sequences of SARS-CoV-2 reported  
329 from December to March 05, 2020 and deposited in the GISAID database. The detailed  
330 analysis of 247 high quality genome sequences of SARS-CoV-2 provides insight into the  
331 evolution and transmission of this novel virus (**Fig 5**). The ancestor of SARS-CoV-2 could  
332 have emerged at a date as early as November, 2019 based on results of phylodynamics analysis  
333 of these genome sequences. According to the time line of outbreaks, the original virus from  
334 Wuhan city and HNSM was responsible for the widespread transmission of SARS-CoV-2 in  
335 different parts of the world in January. The origin of the outbreak was not limited to HNSM,  
336 those which occurred in multiple sites in Wuhan city might have contributed to the early  
337 transmission events and subsequent dissemination to different parts of China and various  
338 countries around the world. These data also implied that wild animals in HNSM may not be  
339 the intermediate host of SARS-CoV-2 as sources other than HNSM are also considered the  
340 origin of this virus. Given the fact that multiple patients in Wuhan were simultaneously infected  
341 by viruses of different genetic composition in the initial outbreak, we hypothesize that a  
342 common wild animal such as wild rat would be the most likely intermediate host. Alternatively,  
343 a common environmental factor, such as a faulty sewage system, may be involved. It is  
344 therefore necessary to investigate the role of a common animal vector or dissemination route  
345 in eliciting the SARS-CoV-2 outbreak.

346

347 Interestingly, as the original virus continued to transmit in China and over the world, it has  
348 evolved into four major genetic clusters, namely super-transmitter clusters, along with other  
349 non-cluster variants derived from the original virus. Each ST cluster carried its unique signature  
350 mutation(s), which enable us to trace its origin and transmission dynamics. In the early  
351 transmission stage (December and January, 2020), variants from the original virus were

352 dominant, yet by the end of February and early March, these four super-transmitters became  
353 dominant, with different STs being prevalent in different regions of the world. ST1 was  
354 prevalent in China and other parts of Asia and became the major virus that caused severe  
355 outbreaks in Washington and California states in the US and South Korea; ST2 and ST3 were  
356 extensively transmitted in other parts of Asia and Europe during the end of January and early  
357 February but its prevalence dropped at the end of February and early March, and was replaced  
358 by ST4 which contributed to the pandemic in Europe. Mapping the mutational profile of viral  
359 genomes enables us to trace the transmission of viruses of different clusters in different parts  
360 of the world. Interestingly, ST4 was not reported in China or other parts of the world. The first  
361 genome of this cluster was reported in Germany and contributed to the rapid dissemination of  
362 SARS-CoV-2 in Europe. Mutation profile with ST4 is unique, with three mutations being  
363 observed in the first viral genome. Genomes with one or two of these three mutations were not  
364 reported anywhere. These data do not simply imply that ST4 originated from Europe. One  
365 limitation of the study is that we can only utilize currently available genome sequences. The  
366 lack of genome sequence of ST4 in other continent does not necessarily mean that ST4 viruses  
367 are not present in other continents. A second limitation of this study is the lack of data to explain  
368 the mechanisms underlying the evolution of these genetic clusters into super-transmitters.  
369 Every ST cluster carried at least one amino acid mutation in different protein. Whether such  
370 mutational changes were the key step that enabled SARS-CoV-2 to evolve into super-  
371 transmitters must be investigated in future research studies.

372

373 Lastly, our data also provided insight into the major transmitting viruses in current pandemic  
374 areas in the world. For example, in Italy, ST2, ST3 were reported in the end of January, while  
375 ST4 was reported in February and early March. Similar trends were seen in other countries  
376 with exception that a higher proportion of the original viral genomes were reported in

377 Netherland. In the US, the original viruses were reported in other states, while ST1 was the  
378 major virus that caused outbreak in Washington and California States. Other ST genomes were  
379 also sporadically reported in the US. Although data from Iran is not available, two genomes  
380 reported from Australia with travel history from Iran were shown to belong to ST3, suggesting  
381 that this cluster was responsible for the pandemic in Iran. In Australia, all genomes were  
382 reported except ST4.

383

384 Using the signature mutations as markers for different STs, we were able to analysed 1539  
385 genomes reported in March. The data further confirmed that four STs became dominant in  
386 March with around 90% of the genomes belonging to these four STs, among which ST4  
387 became the dominant cluster transmitting in Europe. It also started transmitting to other parts  
388 of world including Africa, Asia, North America and Oceania. ST1 was still the major type  
389 transmitting in the US and has transmitted to South America in particular Brazil. These data  
390 confirmed that ST1 and ST4 would become worldwide transmitter and dominate the future  
391 transmission of SARS-CoV-2 in the world.

392

393 In conclusion, this study show that four major genetic clusters of viruses evolved from the  
394 original SARS-CoV-2 and have transmitted extensively over the world, each becoming  
395 dominant in different parts of the world, and that viruses without any signature mutation of the  
396 four super-transmitters appear to be transmitted much less efficiently. These super-transmitters  
397 exhibit not only high transmission efficiency, but also high mutation rate without  
398 compromising infectivity, compromising effectiveness of current infection control effort. Our  
399 findings therefore provide important insight into the molecular features of the highly  
400 transmissible variants of SARS-CoV-2.

401

## 402 References

- 403 1. Cui J, Li F, Shi Z-L: **Origin and evolution of pathogenic coronaviruses.** *Nature reviews*  
404 *Microbiology* 2019, **17**(3):181-192.
- 405 2. Gorbalenya AE: **Severe acute respiratory syndrome-related coronavirus–The species and its**  
406 **viruses, a statement of the Coronavirus Study Group.** *bioRxiv* 2020.
- 407 3. Guan W-j, Ni Z-y, Hu Y, Liang W-h, Ou C-q, He J-x, Liu L, Shan H, Lei C-l, Hui DS: **Clinical**  
408 **characteristics of 2019 novel coronavirus infection in China.** *medRxiv* 2020.
- 409 4. Sardar R, Satish D, Birla S, Gupta D: **Comparative analyses of SAR-CoV2 genomes from**  
410 **different geographical locations and other coronavirus family genomes reveals unique**  
411 **features potentially consequential to host-virus interaction and pathogenesis.** *bioRxiv* 2020.
- 412 5. Lu R, Zhao X, Li J, Niu P, Yang B, Wu H, Wang W, Song H, Huang B, Zhu N: **Genomic**  
413 **characterisation and epidemiology of 2019 novel coronavirus: implications for virus origins**  
414 **and receptor binding.** *The Lancet* 2020.
- 415 6. Zhu N, Zhang D, Wang W, Li X, Yang B, Song J, Zhao X, Huang B, Shi W, Lu R: **A novel**  
416 **coronavirus from patients with pneumonia in China, 2019.** *New England Journal of Medicine*  
417 2020.
- 418 7. Chen N, Zhou M, Dong X, Qu J, Gong F, Han Y, Qiu Y, Wang J, Liu Y, Wei Y *et al*: **Epidemiological**  
419 **and clinical characteristics of 99 cases of 2019 novel coronavirus pneumonia in Wuhan,**  
420 **China: a descriptive study.** *Lancet* 2020.
- 421 8. Shu Y, McCauley J: **GISAID: Global initiative on sharing all influenza data—from vision to**  
422 **reality.** *Eurosurveillance* 2017, **22**(13).
- 423 9. Zhou P, Yang XL, Wang XG, Hu B, Zhang L, Zhang W, Si HR, Zhu Y, Li B, Huang CL *et al*: **A**  
424 **pneumonia outbreak associated with a new coronavirus of probable bat origin.** *Nature* 2020.
- 425 10. X Tang CW, X Li, Y Song, X Yao, X Wu, Y Duan, H Zhang, Y Wang, Z Qian, J Cui, J Lu: **On the**  
426 **origin and continuing evolution of SARS-CoV-2.** *National Science Review*, *nwaa036*,  
427 <https://doi.org/10.1093/nsr/nwaa036> 2020.
- 428 11. Tcherepanov V, Ehlers A, Upton C: **Genome Annotation Transfer Utility (GATU): rapid**  
429 **annotation of viral genomes using a closely related reference genome.** *BMC Genomics* 2006,  
430 **7**:150.
- 431 12. Wu F, Zhao S, Yu B, Chen YM, Wang W, Song ZG, Hu Y, Tao ZW, Tian JH, Pei YY *et al*: **A new**  
432 **coronavirus associated with human respiratory disease in China.** *Nature* 2020.
- 433 13. Hadfield J, Megill C, Bell SM, Huddleston J, Potter B, Callender C, Sagulenko P, Bedford T,  
434 Neher RA: **Nextstrain: real-time tracking of pathogen evolution.** *Bioinformatics* 2018,  
435 **34**(23):4121-4123.
- 436 14. Katoh K, Asimenos G, Toh H: **Multiple alignment of DNA sequences with MAFFT.** In:  
437 *Bioinformatics for DNA sequence analysis.* edn.: Springer; 2009: 39-64.
- 438 15. Nguyen L-T, Schmidt HA, Von Haeseler A, Minh BQ: **IQ-TREE: a fast and effective stochastic**  
439 **algorithm for estimating maximum-likelihood phylogenies.** *Molecular biology and evolution*  
440 2015, **32**(1):268-274.
- 441 16. Sagulenko P, Puller V, Neher RA: **TreeTime: Maximum-likelihood phylodynamic analysis.**  
442 *Virus evolution* 2018, **4**(1):vex042.
- 443 17. Bah T: **Inkscape: guide to a vector drawing program:** prentice hall press; 2007.
- 444 18. Stamatakis A: **RAxML version 8: a tool for phylogenetic analysis and post-analysis of large**  
445 **phylogenies.** *Bioinformatics* 2014, **30**(9):1312-1313.
- 446 19. Letunic I, Bork P: **Interactive tree of life (iTOL) v3: an online tool for the display and**  
447 **annotation of phylogenetic and other trees.** *Nucleic Acids Research* 2016, **44**(Web Server  
448 issue):W242-W245.

449

450

451 **Acknowledgments**

452 We are grateful for the discussion and comments from Prof. Mengsu Yang from City  
453 University of Hong Kong. We acknowledge the use of the genome sequences in the GISAID  
454 platform. This study was not supported by any research grant.

455

456 **Conflicts of interest**

457 We declare that we have no conflict of interest.

458

459 **Table 1. Mutational profile analysis of SARS-CoV-2 genome sequences obtained in**  
 460 **December of 2019.**

Accession ID	Location	Collection date	Host	Origin	Number of mutations / % of sequence homology with reference sequence	Nucleotide changes
EPI_ISL_402119	China / Hubei / Wuhan	2019-12-30	Human	HNSM	0 / 100%	/
EPI_ISL_402124	China / Hubei / Wuhan	2019-12-30	Human	HNSM	0 / 100%	/
EPI_ISL_402125	China	2019-12-31	Human	HNSM	0 / 100%	/
EPI_ISL_402129	China / Hubei / Wuhan	2019-12-30	Human	HNSM	0 / 100%	/
EPI_ISL_403929	China / Hubei / Wuhan	2019-12-30	Human	WH	0 / 100%	/
EPI_ISL_412899	China / Hubei / Wuhan	2019-12-30	human	WH	0 / 100%	/
EPI_ISL_402132	China / Hubei / Wuhan	2019-12-30	Human	WH	1 / 99%	T21656A
EPI_ISL_403930	China / Hubei / Wuhan	2019-12-30	Human	WH	1 / 99%	T6996C
EPI_ISL_412898	China / Hubei / Wuhan	2019-12-30	human	WH	1 / 99%	A24325G
EPI_ISL_402121	China / Hubei / Wuhan	2019-12-30	Human	HNSM	2 / 99%	G20670A; G20679A
EPI_ISL_402127	China / Hubei / Wuhan	2019-12-30	Human	HNSM	2 / 99%	G21316A; A24325G
EPI_ISL_402128	China / Hubei / Wuhan	2019-12-30	Human	HNSM	2 / 99%	G7016A; A21137G
EPI_ISL_402130	China / Hubei / Wuhan	2019-12-30	Human	HNSM	2 / 99%	A8001C; C9534T
EPI_ISL_406798	China / Hubei / Wuhan	2019-12-26	Human	WH	2 / 99%	C6968A; T11764A
EPI_ISL_402123	China / Hubei / Wuhan	2019-12-24	Human	WH	3 / 99%	A3778G; A8388G; T8987A
EPI_ISL_403931	China / Hubei / Wuhan	2019-12-30	Human	WH	6 / 99%	T104A; T111C; T112G; C119G; T120C; G124A

461

462 HNSM, Hua Nan Seafood Wholesale Market.

463

**Table 2. Mutational analysis of genome sequences in super-transmitter cluster 1.**

Accession ID	Location	Collection date	Origin	Number of mutations / % of sequence homology with reference sequence	Nucleotide changes
EPI_ISL_406801	Asia / China / Hubei / Wuhan	2020-01-05	NC	2 / 99%	☆ (C8782T, T28144C)
EPI_ISL_407893	Oceania / Australia / New South Wales / Sydney	2020-01-24	NC	2 / 99%	☆
EPI_ISL_412979	China / Hubei / Wuhan	2020-01-18	NC	2 / 99%	☆
EPI_ISL_413691	China	2020-01	NC	2 / 99%	☆
EPI_ISL_413729	China	2020-02	NC	2 / 99%	☆
EPI_ISL_413746	China	2020-02	NC	2 / 99%	☆
EPI_ISL_413748	China	2020-02	NC	2 / 99%	☆
EPI_ISL_413750	China	2020-02	NC	2 / 99%	☆
EPI_ISL_403932	China / Guandong / Shenzhen	2020-01-14	NC	3 / 99%	☆ ¶ (C8782T, T28144C, C29095T)
EPI_ISL_403933	China / Guandong / Shenzhen	2020-01-15	NC	3 / 99%	☆ ¶
EPI_ISL_403935	China / Guangdong / Shenzhen	2020-01-15	NC	3 / 99%	☆ ¶
EPI_ISL_406030	China / Guangdong / Shenzhen	2020-01-10	WH	3 / 99%	☆ ¶
EPI_ISL_406593	Asia / China / Guandong / Shenzhen	2020-01-13	NC	3 / 99%	☆ ¶
EPI_ISL_406223	USA / Arizona / Phoenix	2020-01-22	NC	4 / 99%	☆ ¶ G11083T
EPI_ISL_413809	China	2020-02	NC	4 / 99%	☆ ¶ G18686T
EPI_ISL_408666	Japan / Tokyo	2020-01-31	NC	4 / 99%	☆ ¶ C2662T
EPI_ISL_408665	Japan / Tokyo	2020-01-29	NC	5 / 99%	☆ ¶ C2662T; C3792T
EPI_ISL_408667	Japan / Tokyo	2020-01-31	NC	5 / 99%	☆ ¶ C2662T, G29705T
EPI_ISL_405839	China / Guangdong / Shenzhen	2020-01-11	WH	5 / 99%	☆ ¶ C9561T; T15607C
EPI_ISL_411956	North America / USA / Texas	2020-02-11	NC	7 / 99%	☆ ¶ T18603C; T18975A; A19175C; C27925T

EPI_ISL_404895	USA / Washington / Snohomish County	2020-01-19	NC	3 / 99%	☆ T28144C
EPI_ISL_407976	Europe / Belgium / Leuven	2020-02-03	WH	3 / 99%	☆ A29863T
EPI_ISL_408480	China / Yunnan / Kunming	2020-01-17	NC	3 / 99%	☆ G11083T
EPI_ISL_408489	Taiwan / Taipei	2020-01-31	WH	3 / 99%	☆ G11528S
EPI_ISL_410535	Singapore	2020-02-03	NC	3 / 99%	☆ G28878A
EPI_ISL_411926	Taiwan / Taipei	2020-01-24	NC	3 / 99%	☆ A29889G
EPI_ISL_413854	China / Guangdong	2020-01-30	NC	3 / 99%	☆ C6501T
EPI_ISL_411060	China / Fujian	2020-01-21	NC	3 / 99%	☆ ⊥ (C8782T, T28144C, C18060T)
EPI_ISL_407214	USA / Washington	2020-01-25	NC	3 / 99%	☆ ⊥
EPI_ISL_407215	USA / Washington	2020-01-25	NC	3 / 99%	☆ ⊥
EPI_ISL_408478	China / Chongqing / Yongchuan	2020-01-21	NC	5 / 99%	☆ ⊥ C29200T, C1342T
EPI_ISL_413456	USA / Washington / King County	2020-02-20	NC	5 / 99%	☆ ⊥ ⊚ (C8782T, T28144C, C18060T C17747T, A17858G)
EPI_ISL_413560	USA / Washington	2020-02-28	NC	5 / 99%	☆ ⊥ ⊚
EPI_ISL_412970	USA / Washington / Snohomish County	2020-02-24	NC	6 / 99%	☆ ⊥ ⊚ C5784T,
EPI_ISL_413457	USA / Washington	2020-02-29	NC	6 / 99%	☆ ⊥ ⊚ C20S,
EPI_ISL_413458	USA / Washington	2020-03-01	NC	6 / 99%	☆ ⊥ ⊚ T20281C
EPI_ISL_413563	USA / Washington	2020-03-03	NC	6 / 99%	☆ ⊥ ⊚ C9430A,
EPI_ISL_413650	USA / Washington	2020-03-05	NC	6 / 99%	☆ ⊥ ⊚ T23010C
EPI_ISL_413651	USA / Washington	2020-03-05	NC	6 / 99%	☆ ⊥ ⊚ T23010C
EPI_ISL_413653	USA / Washington	2020-03-05	NC	6 / 99%	☆ ⊥ ⊚ A6T
EPI_ISL_413455	USA / Washington	2020-02-28	NC	8 / 99%	☆ ⊥ ⊚ T29867A, G29868A, C29870A
EPI_ISL_413486	USA / Washington	2020-03-01	NC	8 / 99%	☆ ⊥ ⊚ A3406C, C5784T, C23525T,

EPI_ISL_413925	USA / California / San Francisco	2020-03-05	GPCS	8 / 99%	☆ ≡ ⊥ C23185T, A3046G, A16467G,
EPI_ISL_413931	USA / California / San Francisco	2020-03-05	NC	9 / 99%	☆ ≡ ⊥ A3046G, A16467G, G16975T, C23185T,
EPI_ISL_413562	USA / Washington	2020-03-02	NC	10 / 99%	☆ ≡ ⊥ C313del, C9180T, G29864A, T29867A, G29868A
EPI_ISL_413652	USA / Washington	2020-03-05	NC	11 / 99%	☆ ≡ ⊥ T23010C, G29861A, G29864C, T29867A, G29868C, C29870A
EPI_ISL_407193	South Korea / Gyeonggi-do	2020-01-25	NC	4 / 99%	☆ ⊥ (C8782T, T28144C, T4402C; G5062T)
EPI_ISL_412870	South Korea/ Seoul	2020-01-30	NC	4 / 99%	☆ ⊥
EPI_ISL_413513	South Korea	2020-02-27	NC	4 / 99%	☆ ⊥
EPI_ISL_413514	South Korea	2020-02-27	NC	4 / 99%	☆ ⊥
EPI_ISL_413515	South Korea	2020-02-27	NC	4 / 99%	☆ ⊥
EPI_ISL_413516	South Korea	2020-02-27	NC	4 / 99%	☆ ⊥
EPI_ISL_413518	China / Beijing	2020-01-26	NC	4 / 99%	☆ ⊥
EPI_ISL_413519	China / Beijing	2020-01-28	NC	4 / 99%	☆ ⊥
EPI_ISL_413521	China / Beijing	2020-01-28	NC	4 / 99%	☆ ⊥
EPI_ISL_413520	China / Beijing	2020-01-28	NC	5 / 99%	☆ ⊥ A29301T
EPI_ISL_412871	South Korea / Seoul	2020-01-31	NC	6 / 99%	☆ ⊥ C1779T, C15017T
EPI_ISL_410718	Queensland / Gold Coast	2020-02-05	NC	4 / 99%	☆ ⊥ (C8782T, T28144C, G28878A; G29742A)
EPI_ISL_411954	USA / California	2020-02-06	NC	4 / 99%	☆ ⊥
EPI_ISL_413853	China / Guangdong	2020-01-30	NC	4 / 99%	☆ ⊥
EPI_ISL_410717	Australia / Queensland / Gold Coast	2020-02-05	NC	5 / 99%	☆ ⊥ T28144C
EPI_ISL_412978	China / Hubei / Wuhan	2020-01-17	NC	4 / 99%	☆ C12141A, C23816T
EPI_ISL_413711	China	2020-02	NC	4 / 99%	☆ C6501T, C16887T

EPI_ISL_413523	India / Kerala	2020-01-31	China	6 / 99%	☆ A1691G, C6501T, C16877T, C24351T
EPI_ISL_413749	China	2020-02	NC	4 / 99%	☆ C14768T, A17805T
EPI_ISL_413858	China / Guangdong	2020-01-30	NC	4 / 99%	☆ A27749N, G27750N
EPI_ISL_412980	China / Hubei / Wuhan	2020-01-18	NC	5 / 99%	☆ T18996C, C24370T, T29029C
EPI_ISL_407071	Europe / England	2020-01-29	NC	5 / 99%	☆ T22586Y; T23605G; T28144C
EPI_ISL_407073	Europe / England	2020-01-29	NC	5 / 99%	☆ T23605G; T18488C, A29596G
EPI_ISL_412982	China / Hubei / Wuhan	2020-02-07	NC	5 / 99%	☆ G5657A, A23403G, A25725G,
EPI_ISL_413697	China	2020-02	NC	5 / 99%	☆ C207T, T946C, A11430G
EPI_ISL_413751	China	2020-02	NC	5 / 99%	☆ TTT27792-27794del
EPI_ISL_413761	China	2020-02	NC	5 / 99%	☆ C207T, T946C, A11430G
EPI_ISL_408484	China / Sichuan / Chengdu	2020-01-15	NC	6 / 99%	☆ ■ T19190A; C24034T
EPI_ISL_406034	USA / California / Los Angeles	2020-01-23	NC	7 / 99%	☆ ■ G1548A; C24034, A28792T
EPI_ISL_410045	USA / Illinois	2020-01-28	NC	7 / 99%	☆ ■ T490A; C3177T, C24034T;
EPI_ISL_412028	Hong Kong	2020-01-22	NC	7 / 99%	☆ ■ C1663T, G22661T, G29862T
EPI_ISL_408668	Vietnam / Thanh Hoa	2020-01-24	NC	11 / 99%	☆ ■ A27T; C28del; C24034T; T29858C; G29861C; G29864del; T29867A
EPI_ISL_407896	Queensland / Gold Coast	2020-01-30	NC	7 / 99%	☆ A21949M; C24790T; C25587T; G28878A; G29742A

EPI_ISL_412873	South Korea / Chungcheongnam-do	2020-02-06	NC	7 / 99%	☆ T3086C, C6255T, G11083T, G17122A, A29871G
EPI_ISL_413791	China	2020-02	NC	7 / 99%	☆ C207T, T946C, A11430G, A16474G, C25000A
EPI_ISL_404253	USA / Illinois / Chicago	2020-01-21	NC	8 / 99%	T490W; C3177Y; C24034Y; T26729Y; G28077Y; C28854Y
EPI_ISL_412869	South Korea /Seoul	2020-01-30	NC	8 / 99%	☆ A1740C, C8782T, C17104T, G26167T, G29593A, A29869G
EPI_ISL_412983	China / Hubei / Tianmen	2020-02-08	NC	9 / 99%	☆ A3175G, G3179A, C14422T, C14585T, G23405C, C28315T, T29680K
EPI_ISL_413485	China / Anhui / Suzhou	2020-01-24	NC	9 / 99%	☆ A4T, C2189T, T3086C, A5094G, G11083del, C16049T, G17122A, ☆ A6604R; C13681M; A13682M; C13684M; T13686K; G13687K; A13693W; G28878A; G29742A
EPI_ISL_407894	Queensland / Gold Coast	2020-01-28	NC	11 / 99%	

466

467 ☆, C8782T, T28144C; †, C29095T; ‡ C18060T; § C17747T, A17858G; ¶, T4402C; G5062T; ††, G28878A; G29742A; †††, T26729C;

468 NC, Not confirmed;

**Table 3. Mutational analysis of genome sequences in super-transmitter cluster 2.**

Accession ID	Location	Collection date	Origin	Number of mutations / % of sequence homology with reference sequence	Nucleotide changes
EPI_ISL_408977	Australia / Sydney	2020-01-25	NC	1 / 99%	○ (G26144T)
EPI_ISL_406036	USA / California	2020-01-22	NC	2 / 99%	○ C17000T
EPI_ISL_412029	Hong Kong	2020-01-30	NC	2 / 99%	○ T13929C
EPI_ISL_413863	China / Guangdong	2020-02-01	NC	2 / 99%	○ C22787G
EPI_ISL_406596	France / Paris	2020-01-23	NC	2 / 99%	○ ▽ (G26144T, G22661T)
EPI_ISL_406597	France / Paris	2020-01-23	WH	2 / 99%	○ ▽
EPI_ISL_410720	France / Paris	2020-01-23	NC	2 / 99%	○ ▽
EPI_ISL_411219	France / Paris	2020-01-28	NC	2 / 99%	○ ▽
EPI_ISL_411220	France / Paris	2020-01-28	NC	2 / 99%	○ ▽
EPI_ISL_410713	Singapore	2020-01-27	NC	2 / 99%	○ C28849T
EPI_ISL_410714	Singapore	2020-02-03	NC	2 / 99%	○ C21859T
EPI_ISL_410536	Singapore	2020-02-06	NC	2 / 99%	○ ▼ (G26144T, C21859T)
EPI_ISL_410715	Singapore	2020-02-04	NC	2 / 99%	○ ▼
EPI_ISL_410716	Singapore	2020-02-04	NC	2 / 99%	○ ▼
EPI_ISL_410546	Italy / Rome	2020-01-31	HB	2 / 99%	○ ◇ (G26144T, G11083T)
EPI_ISL_412974	Italy / Rome	2020-01-29	NC	2 / 99%	○ ◇
EPI_ISL_410545	Italy / Rome	2020-01-29	HB	3 / 99%	○ A2269T; G11083N;
EPI_ISL_413603	Finland / Helsinki	2020-03-03	NC	4 / 99%	○ ◇ C14805T, G29405C
EPI_ISL_413016	Brazil / Sao Paulo	2020-02-28	Italy	5 / 99%	○ ◇ C2388T, C14805T, T17247C,
EPI_ISL_413019	Switzerland / Zurich	2020-02-26	NC	9 / 99%	○ ◇, G11084TTTin, C14805T, T17247C, C24378T, C26894T
EPI_ISL_413025	USA / Washington	2020-02-27	NC	12 / 99%	○ ◇, A35T, C36T, T2446C, C3411T, G5572T, C14805T, G29864A, T29867A, G29868A, C29870A

EPI_ISL_406031	Taiwan / Kaohsiung	2020-01-23	NC	4 / 99%	○ G16188T; A25964G; 29877Tin
EPI_ISL_413018	South Korea	2020-02-06	NC	4 / 99%	○ A2707G, G26640T, T26677C
EPI_ISL_412116	England	2020-02-09	NC	5 / 99%	○ A2470G, C2558T, G11083N, C14805T,
EPI_ISL_413017	South Korea	2020-02-06	NC	6 / 99%	○ T4402C, G5062T, G26640T, T26677C, T28144C
EPI_ISL_411951	Sweden	2020-02-07	NC	7 / 99%	○ G2717A; A9274G; C13225G; T13226C; A17376G; T23952G;
EPI_ISL_411929	South Korea	2020-01	WH	9 / 99%	○ G2971T; C6031T; C12115T; T15597C; C20936T; C22224G; G25775T; T26354A
EPI_ISL_406844	Australia / Victoria	2020-01-25	NC	13 / 99%	○ T19065C; T22303G; 29750-29759del

471 ○, G26144T; ▽, G22661T; ▼, C21859T; ◇, G11083T;

472

**Table 4. Mutational analysis of genome sequences in super-transmitter cluster 3.**

Accession ID	Location	Origin	Collection date	Number of mutations / % of sequence homology with reference sequence	Nucleotide changes
EPI_ISL_408481	China / Chongqing	NC	2020-01-18	1 / 99%	⊙ (G11083T)
EPI_ISL_407988	Singapore	NC	2020-02-01	1 / 99%	⊙
EPI_ISL_412968	Japan	NC	2020-02-10	1 / 99%	⊙
EPI_ISL_410546	Italy / Rome	HB	2020-01-31	2 / 99%	⊙ G26144T
EPI_ISL_412030	Hong Kong	NC	2020-02-01	2 / 99%	⊙ G29841A
EPI_ISL_412969	Japan	NC	2020-02-10	2 / 99%	⊙ C29635T
EPI_ISL_412974	Italy / Rome	NC	2020-01-29	2 / 99%	⊙ G26144T
EPI_ISL_408430	France / Paris	NC	2020-01-29	3 / 99%	⊙ ⊙ (G11083T, C1190T, C9438T)
EPI_ISL_410984	France / Paris	NC	2020-01-29	3 / 99%	⊙ ⊙
EPI_ISL_411218	France / Paris	NC	2020-02-02	3 / 99%	⊙ ⊙
EPI_ISL_408480	China / Yunnan / Kunming	NC	2020-01-17	3 / 99%	⊙ † (G11083T, C8782T, T28144C)
EPI_ISL_406223	USA / Arizona / Phoenix	NC	2020-01-22	4 / 99%	⊙ † C29095T
EPI_ISL_412873	South Korea	NC	2020-02-06	7 / 99%	⊙ † T3086C, C6255T, G17122A, A29871G
EPI_ISL_413485	China / Anhui / Suzhou	NC	2020-01-24	9 / 99%	⊙ † A4T, C2189T, T3086C, A5094G, C16049T, G17122A
EPI_ISL_413603	Finland / Helsinki	NC	2020-03-03	4 / 99%	⊙ \$ G26144T, G29405C
EPI_ISL_413016	Brazil / Sao Paulo	Italy	2020-02-28	5 / 99%	⊙ \$ C2388T, T17247C, G26144T
EPI_ISL_413025	USA / Washington	NC	2020-02-27	12 / 99%	⊙ \$ A35T, C36T, T2446C, C3411T, G5572T, G26144T, G29864A, T29867A, G29868A, C29870A
EPI_ISL_413214	Australia Sydney	NC	2020-02-29	5 / 99%	⊙ % G29374A

EPI_ISL_412975	Australia Sydney	Iran	2020-02-28	6 / 99%	⊙ % G4255A, , A20047G
EPI_ISL_413213	Australia / Sydney	Iran	2020-02-29	7 / 99%	⊙ % C884T, G8653T, C24704T
EPI_ISL_408482	Shandong / Qingdao	NC	2020-01-19	7 / 99%	⊙ % C7299T; C27612G; T28688C
EPI_ISL_413589	Netherlands / Utrecht	NC	2020-03-01	8 / 99%	⊙ C241T, G2527T, C3037T, C6428T, C14408T, A23403G, A25575C

---

475

476 ⊙, G11083T; ©, C1190T, C9438T; ð, C8782T, T28144C; \$, C14805T; %, G1397A, T28688C, G29742T.

**Table 5. Mutational analysis of genome sequences in super-transmitter cluster 4.**

Accession ID	Location	Origin	Collection date	Number of mutations / % of sequence homology with reference sequence	Nucleotide changes
EPI_ISL_406862	Germany / Bavaria / Munich	NC	2020-01-28	3 / 99%	Δ (C241T, C3037T, A23403G)
EPI_ISL_413555	United Kingdom / Wales	NC	2020-02-27	4 / 99%	Δ # (C241T, C3037T, A23403G, C14408T)
EPI_ISL_413566	Netherlands / Blaricum	NC	2020-03-02	4 / 99%	Δ #
EPI_ISL_413591	Netherlands / Zeewolde	NC	2020-03-02	4 / 99%	Δ #
EPI_ISL_413593	Luxembourg	NC	2020-02-29	5 / 99%	Δ # C23575T
EPI_ISL_413648	Portugal	Spain	2020-03-01	5 / 99%	Δ # C29144T
EPI_ISL_412973	Italy	NC	2020-02-20	6 / 99%	Δ # T29867N, G29868N
EPI_ISL_413489	Italy / Milan	NC	2020-03-03	8 / 99%	Δ # A187G, A6956C, T29867N, G29868N
EPI_ISL_413602	Finland / Helsinki	NC	2020-03-03	6 / 99%	Δ # G22865T, C29585T
EPI_ISL_413572	Netherlands / Haarlem	NC	2020-03-01	7 / 99%	Δ # T1666C, C3037T, G25563T
EPI_ISL_413589	Netherlands / Utrecht	NC	2020-03-01	8 / 99%	Δ # G2527T, C6428T, G11083T, A25575C
EPI_ISL_413022	Switzerland / Zurich	NC	2020-02-29	7 / 99%	Δ # & (C241T, C3037T, A23403G, C14408T, G28881A, G28882A, G28883C)
EPI_ISL_413579	Netherlands / Nootdorp	NC	2020-03-03	7 / 99%	Δ # &
EPI_ISL_413587	Netherlands / Tilburg	NC	2020-03-03	7 / 99%	Δ # &
EPI_ISL_412912	Germany/Baden-Wuerttemberg	Italy	2020-02-25	8 / 99%	Δ # & G10265A
EPI_ISL_413584	Netherlands / Rotterdam	NC	2020-03-03	8 / 99%	Δ # & C27046T
EPI_ISL_413647	Portugal	Germany	2020-03-01	8 / 99%	Δ # & C27046T
EPI_ISL_413604	Finland / Helsinki	NC	2020-03-03	9 / 99%	Δ # & C27046T, T29807C
EPI_ISL_413023	Switzerland / Zurich	NC	2020-02-29	9 / 99%	Δ # & A22168C,

EPI_ISL_412972	Mexico / Mexico City	Italy	2020-02-27	10 / 99%	△ # & C13206G, A15807del, G24268del,
----------------	----------------------	-------	------------	----------	---

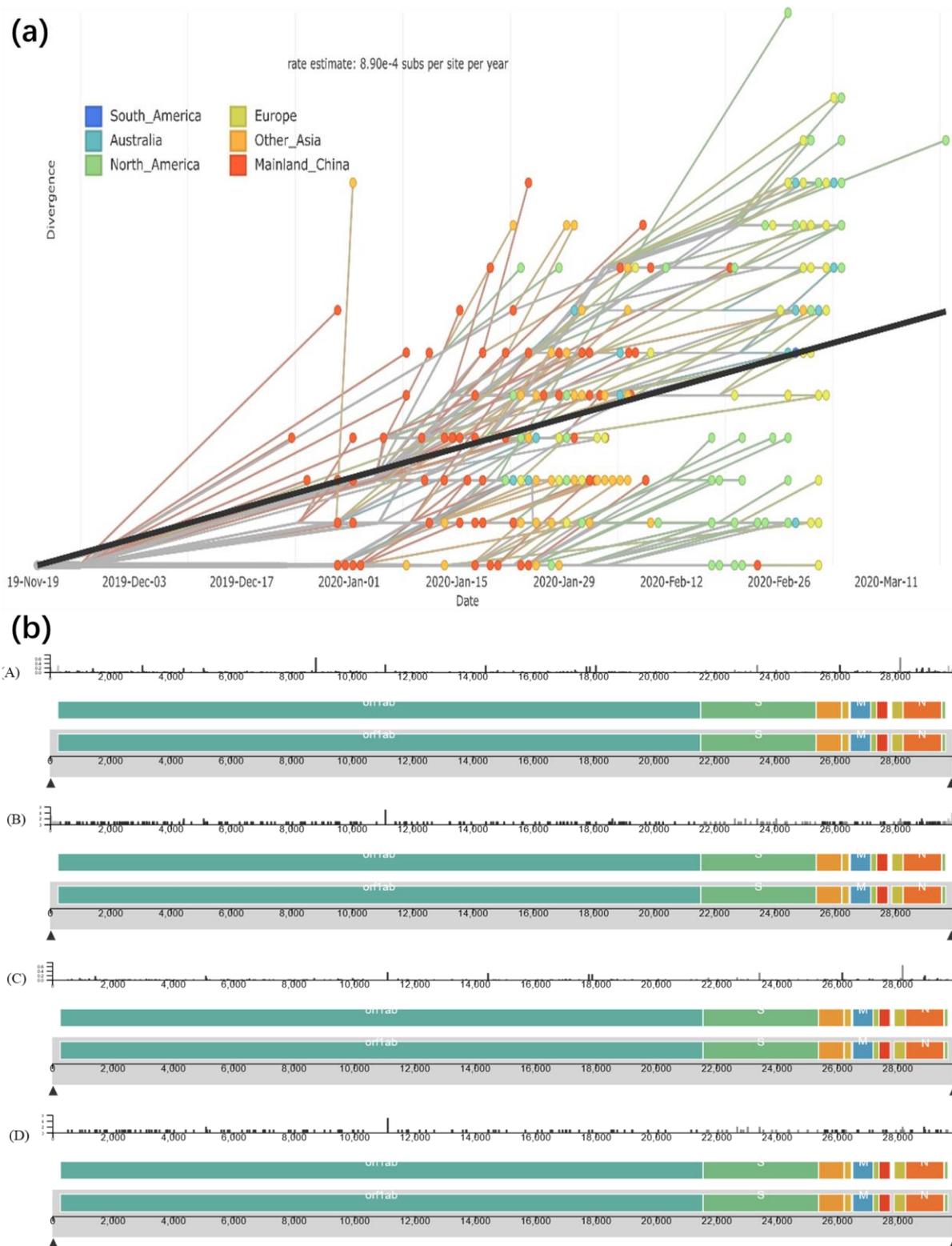
---

479  
480 △, C241T, C3037T, A23403G; #, C14408T; &, C14408T, G28881A, G28882A, G28883C.  
481

482 **Table 6. Distribution of different types of SARS-CoV-2 over the world after February 29, 2020.**

Types of viruses	Numbers of genomes (%)	Major areas	Locations / numbers of genomes (%)						
			Africa	Asia	Europe	NA (USA)	NA (Canada)	Oceania	South America
<b>Total</b>	1539		20(1)	33(2)	968(63)	441(29)	11(1)	27(2)	23(2)
<b>Original derivatives</b>	173 (11)	Europe (UK and Netherland)	0(0)	0(0)	158(16)	11(3)	2(2)	0(0)	2(9)
<b>ST1</b>	340 (21)	USA, Europe, South America	1(5)	1(3)	30(3)	282(64)	0(0)	5(19)	8(35)
<b>ST2</b>	172 (11)	Europe	1(5)	0(0)	111(12)	9(2)	2(2)	3(11)	0(0)
<b>ST3</b>	132 (9)	Asia, Europe	0(0)	17(52)	119(12)	12(3)	5(5)	15(56)	0(0)
<b>ST4</b>	856 (56)	Africa, Asia, Europe, Oceania	18(90)	16(48)	663(69)	137(14)	4(4)	5(19)	12(52)

483



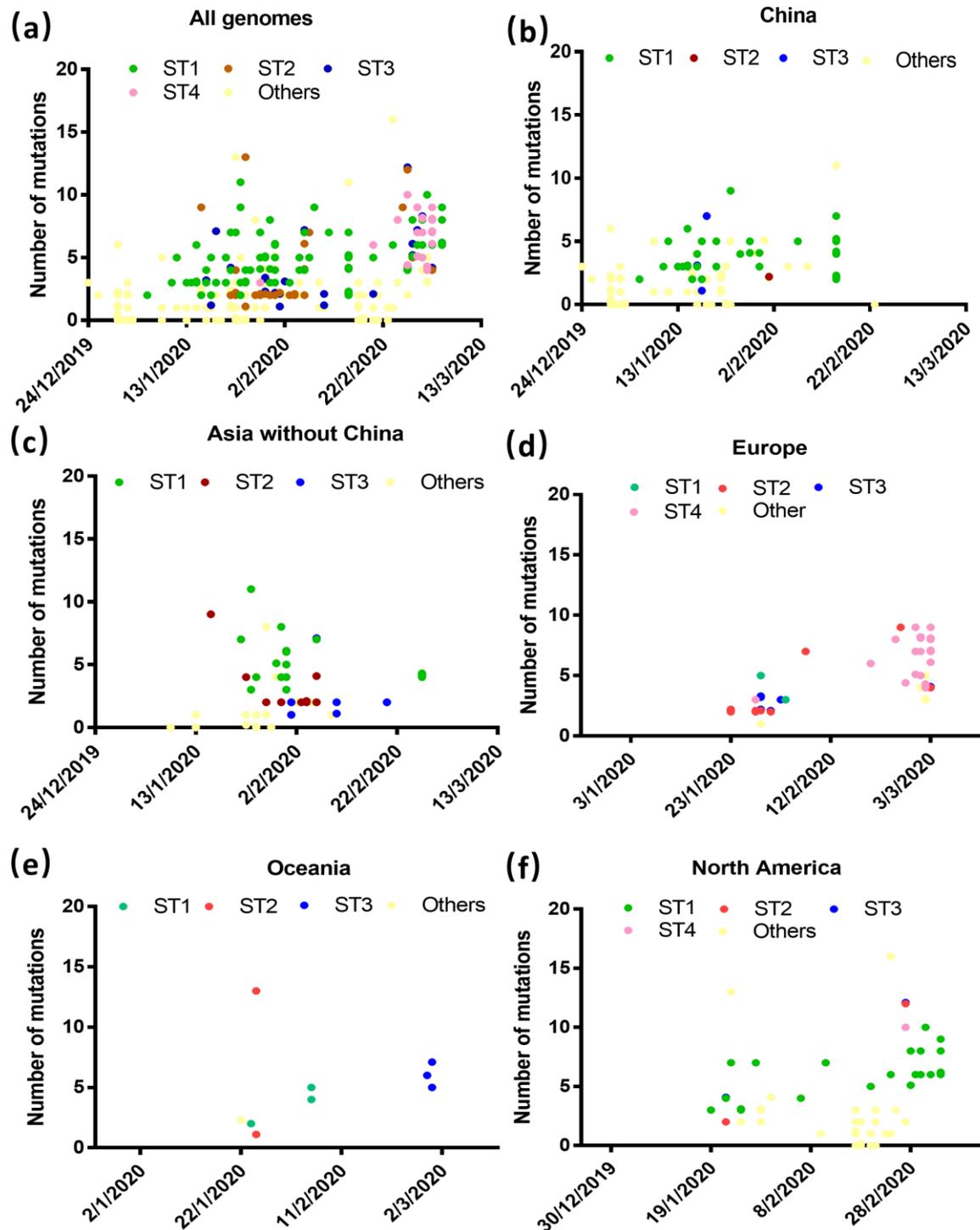
484

485 **Figure 1. Root-to-tip regression scatter plots and distribution of mutations across the**  
486 **SAR-CoV-2 genomes. (a) Root-to-tip regression scatter plots of different strains of SARS-**  
487 **CoV-2. (b) Distribution of mutations across the SAR-CoV-2 genome. (A) and (B) illustrated**  
488 **the entropy and events of nucleotide mutations. (C) and (D) illustrated the entropy and events**  
489 **of amino acid mutations. Figure 1. Root-to-tip regression scatter plots of different strains of**  
490 **SARS-CoV-2. Dots in the plot indicate the SARS-CoV-2 isolates used in this study. The color**

491 of each dot represents the region of isolation of the corresponding isolate. Ancestral state  
492 reconstruction and branch length timing were performed with TreeTime [16].

493





498  
499

500 **Figure 3. Changes in the distribution pattern and mutation rate of different super-**  
501 **transmitter clusters in various continents over time.** Distribution of different STs and their  
502 mutations (a) Overall, (b) in China, (c), in Asian countries excluding China, (d) in Europe, (e)  
503 in Oceania, and (f) in North America. Two genomes with over 20 mutations were not included  
504 to facilitate easy visualization of the graphs.



505

506 **Figure 4. Transmission of super-transmitters and other derivatives of the original SARS-CoV-2 in different areas of world.** The derivatives  
 507 of the original virus have been transmitted worldwide and contributed to the early outbreak of COVID-19. ST1 transmits mainly in Asia and the  
 508 US but was less prevalent in other parts of the world. ST2 and ST3 was transmitted mainly in Asian countries other than China, as well as Europe  
 509 from mid of January to mid of February. ST4 was transmits mainly in Europe in the beginning and then transmitted to all over the world.