



## Abstract

42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62

COVID-19 has widely spread across the world, and much research is being conducted on the causative virus SARS-CoV-2. To help control the infection, we developed the Coronavirus GenBrowser (CGB) to monitor the pandemic. CGB allows visualization and analysis of the latest viral genomic data. Distributed genome alignments and an evolutionary tree built on the existing subtree are implemented for easy and frequent updates. The tree-based data are compressed at a ratio of 2,760:1, enabling fast access and analysis of SARS-CoV-2 variants. CGB can effectively detect adaptive evolution of specific alleles, such as D614G of the spike protein, in their early stage of spreading. By lineage tracing, the most recent common ancestor, dated in early March 2020, of nine strains collected from six different regions in three continents was found to cause the outbreak in Xinfadi, Beijing, China in June 2020. CGB also revealed that the first COVID-19 outbreak in Washington State was caused by multiple introductions of SARS-CoV-2. To encourage data sharing, CGB credits the person who first discovers any SARS-CoV-2 variant. As CGB is developed with eight different languages, it allows the general public in many regions of the world to easily access pre-analyzed results of more than 132,000 SARS-CoV-2 genomes. CGB is an efficient platform to monitor adaptive evolution and transmission of SARS-CoV-2.

## 63 Main Text

64 Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2)<sup>1-3</sup> has infected more  
65 than 75 million people, and at least 1.6 million people in more than 200 countries  
66 have died from COVID-19. Many factors have contributed to the COVID-19  
67 pandemic<sup>4-6</sup>, and it has been predicted that the COVID-19 pandemic may last until  
68 2025<sup>7,8</sup>. The pathogen genomics platform Nextstrain has allowed analysis of genomic  
69 sequences of approximately 4,000 strains of SARS-CoV-2 and investigation of its  
70 evolution<sup>9</sup>. As more than 210,000 SARS-CoV-2 strains have been sequenced (Figure  
71 S1), analysis of all strains has far exceeded the capacity of Nextstrain. New  
72 approaches are needed to accomplish this task.

73  
74 To allow timely analysis of a large number of viral genomes, we first solved the  
75 problem that all viral genomes have to be re-aligned when nucleotide sequences of  
76 new genomes become available. This is extremely time-consuming. With the  
77 distributed alignment system (Figure 1), we dramatically reduced the total time  
78 required for the alignment. We also built the evolutionary tree on the existing tree and  
79 new genomic data in order to reduce the complexity of tree construction. With these  
80 modifications, hundreds of thousands of SARS-CoV-2 genomes can be timely  
81 analyzed with data easily shared and visualized on personal computers and smart  
82 phones (Figure 1).

83  
84 For genomic sequence alignments, 132,443 high quality SARS-CoV-2 genomic  
85 sequences were obtained from the 2019nCoV database<sup>10</sup>, which is an integrated  
86 resource based on CNGBdb, GenBank, GISAID<sup>11,12</sup>, GWH<sup>13</sup>, and NMDC. The  
87 sequences were aligned<sup>14</sup> to that of the reference genome and presented as distributed  
88 alignments. Genomic sequences of the bat coronavirus RaTG13<sup>15</sup>, the pangolin  
89 coronavirus PCoV-GX-PIE<sup>16</sup>, and the early SARS-CoV-2 strains collected before Jan  
90 31, 2020 were jointly used to identify ancestral alleles of SARS-CoV-2. The  
91 evolutionary tree was rebuilt based on new data and the existing tree, and mutations in  
92 strains of each branch were indicated according to the principle of parsimony<sup>17</sup>.

93  
94 The genome-wide mutation rate of coronaviruses has been determined to be  
95  $10^{-4} - 10^{-2}$  per nucleotide per year<sup>18</sup>. As this range of mutation rate is too wide,  
96 we decided to estimate more precisely the genome-wide mutation rate ( $\mu$ ) of  
97 SARS-CoV-2 in a timely manner and determined that  $\mu = 6.753 \times 10^{-4}$  per  
98 nucleotide per year (95% confidence interval:  $4.581 \times 10^{-4}$  to  $9.253 \times 10^{-4}$ ). This  
99 calculation did not require information on demography and the time of appearance of  
100 the most recent common ancestor (MRCA) of SARS-CoV-2. The estimated  $\mu$  was  
101 lower than that of other coronaviruses, such as SARS-CoV (0.80 to  $2.38 \times 10^{-3}$   
102 per nucleotide per year)<sup>18</sup> and MERS-CoV ( $1.12 \times 10^{-3}$  per nucleotide per year)<sup>19</sup>.  
103 It was also lower than that determined by other investigators (1.19 to  $1.31 \times 10^{-3}$   
104 per nucleotide per year)<sup>20</sup>. Various mutation rates were found in different regions of  
105 SARS-CoV-2 genome. The mutation rate of each gene is presented in Table S1.

106

107 Similar to Nextstrain<sup>9</sup> and the WashU Virus Genome Browser<sup>21</sup>, the pre-analyzed  
108 genomic variant data on CGB are shared with the general public. The size of  
109 distributed alignments is 3,894 Mb for the high-quality 132,443 SARS-CoV-2  
110 genomic sequences collected globally. The tree-based data format allows the  
111 compression ratio to reach 2,760:1, meaning that the size of compressed data file is as  
112 small as 1.41 Mb. This approach ensures low-latency access to the data and enables  
113 fast sharing and re-analysis of a large number of SARS-CoV-2 genomic variants.  
114

115 To visualize, search, and filter the results of genomic analysis, both desktop  
116 standalone and web-based user-interface of CGB were developed. Similar to the  
117 UCSC SARS-CoV-2 Genome Browser<sup>22</sup> and the WashU Virus Genome Browser<sup>21</sup>,  
118 six genomic-coordinate annotated tracks were developed to show genome structure  
119 and key domains, allele frequencies, sequence similarity, multi-coronavirus genome  
120 alignment, and primer sets for detection of various SARS-CoV-2 strains. To  
121 efficiently visualize the results of genomic analysis, movie-making ability was  
122 implemented for painting the evolutionary tree, and only elements shown on the  
123 screen and visible to the user would be painted. This design makes the visualization  
124 process highly efficient, and the tree of more than 132,000 viral strains can be  
125 visualized even on a smart phone.  
126

127 CGB detects on-going positive selection based on frequency trajectory of a selected  
128 allele. It has been shown that the spike protein G614 variant has a fitness advantage<sup>23</sup>.  
129 Our analysis using CGB confirmed this finding even when the frequency of this  
130 mutation was very low (< 10%). Moreover, two previously identified variants  
131 (ORF1b:P314L, and N:A220V)<sup>24</sup> and five potentially advantageous variants were  
132 also identified even though their frequency was lower than 10% (Figure 2, Table S3).  
133 Thus, CGB is an efficient monitoring platform for detecting advantageous variants  
134 before they become widely spread (Figure S12).  
135

136 CGB is also an efficient platform to investigate local and global transmission of  
137 COVID-19 (Figure 3). There was a recent outbreak in Qingdao, China<sup>25</sup> after two  
138 dock workers were found to have asymptomatic infections on September 24, 2020.  
139 CGB lineage tracing revealed that the sequence of a sample collected from the outer  
140 packaging of cold-chain products is identical to that of the most recent common  
141 ancestor of the two viral strains isolated from the two dock workers (Figure 3B),  
142 suggesting that infection of these two individuals was cold-chain related. However,  
143 this possibility remains to be determined.  
144

145 CGB lineage tracing also revealed the difficulty in the control of COVID-19  
146 pandemic. There was a recent outbreak in Xinfadi, Beijing, China<sup>26</sup>. The sequences  
147 of two viral isolates (Beijing/IVDC-02-06, Beijing/BJ0617-01-Y), collected from two  
148 Xinfadi cases on June 11 and 14, 2020, were found to be identical to the sequence of  
149 an ancestral strain (Figure 3C) dated on March 6, 2020 (95% CI: February 28 – March  
150 17, 2020). This ancestral strain was found to spread to Taiwan, India, Czech Republic,

151 England, Denmark, and Colombia and caused the outbreak in Beijing three months  
152 later. These two Xinfadi strains were also found to evolve significantly slowly  
153 ( $P = 0.0043$  and  $0.0051$ , respectively) because no mutations were detected during  
154 the three months.

155

156 CGB is a powerful tool for the identification of global and regional routes of virus  
157 transmission as it is specially designed to determine whether the mutation rate of a  
158 specific strain is lower than the average mutation rate of the entire set of strains. This  
159 lineage-specific reduced mutation rate could be due to a long period of dormancy  
160 caused by the yet to be confirmed cold-chain preservation<sup>27</sup> or other reasons. Among  
161 the 132,443 SARS-CoV-2 strains, 4,597 strains were found to evolve significantly  
162 slowly ( $P = 2.18 \times 10^{-8} \sim 0.0041$ , Supplemental excel file) and did not mutate  
163 within at least 100 days. This data showed that CGB can narrow the time period for  
164 tracing the transmission of a specific strain.

165

166 A study on the sequences of 453 SARS-CoV-2 genomes collected before mid-March  
167 2020 suggested that the first COVID-19 outbreak in Washington State was due to a  
168 single introduction<sup>28</sup>. However, results of CGB analysis suggest that the first  
169 Washington State outbreak was actually caused by multiple introductions (Figure  
170 S14).

171

172 All the timely-updated data are freely available at <https://bigd.big.ac.cn/ncov/apis/>.  
173 The free desktop standalone version provides the full function of CGB and has a  
174 plug-in module for the eGPS software (<http://www.egps-software.net/>)<sup>29</sup>. Although  
175 the web-based tool is a simplified version of CGB (Figure 4)  
176 (<https://www.biosino.org/genbrowser/> and <https://bigd.big.ac.cn/genbrowser/>), it  
177 provides a convenient way to access the data via a web browser, such as Google  
178 Chrome, Firefox and Safari. The web-based CGB package can be downloaded and  
179 reinstalled on any websites. For educational purpose, eight language versions  
180 (Chinese, English, German, French, Italian, Portuguese, Russian, and Spanish) are  
181 available.

182

### 183 **Acknowledgments**

184

185 We thank Ya-Ping Zhang for providing valuable advices and encouragement, and the  
186 researchers who generated and deposited the sequencing data of SARS-CoV-2 in  
187 GISAID, GenBank, CNGBdb, GWH, and NMDC, making this study possible. This  
188 work was supported by a grant from the National Key Research and Development  
189 Project (No. 2020YFC0847000).

190

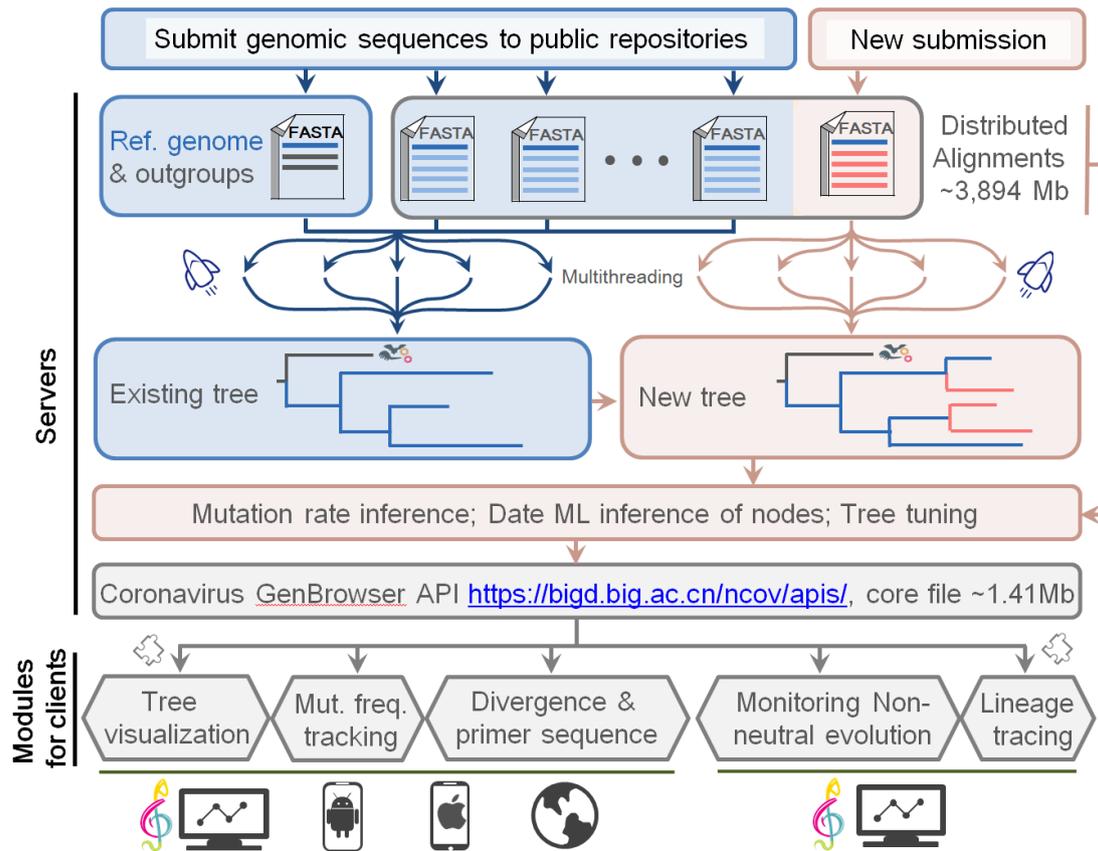
### 191 **Members of the language translation team**

192 German: Ning He<sup>7</sup>, Jing Lv<sup>7</sup>, Ting Peng<sup>7</sup>

193 Italian: Ting Zhou<sup>7</sup>, Nan Yang<sup>7</sup>, Siyi Hou<sup>7</sup>

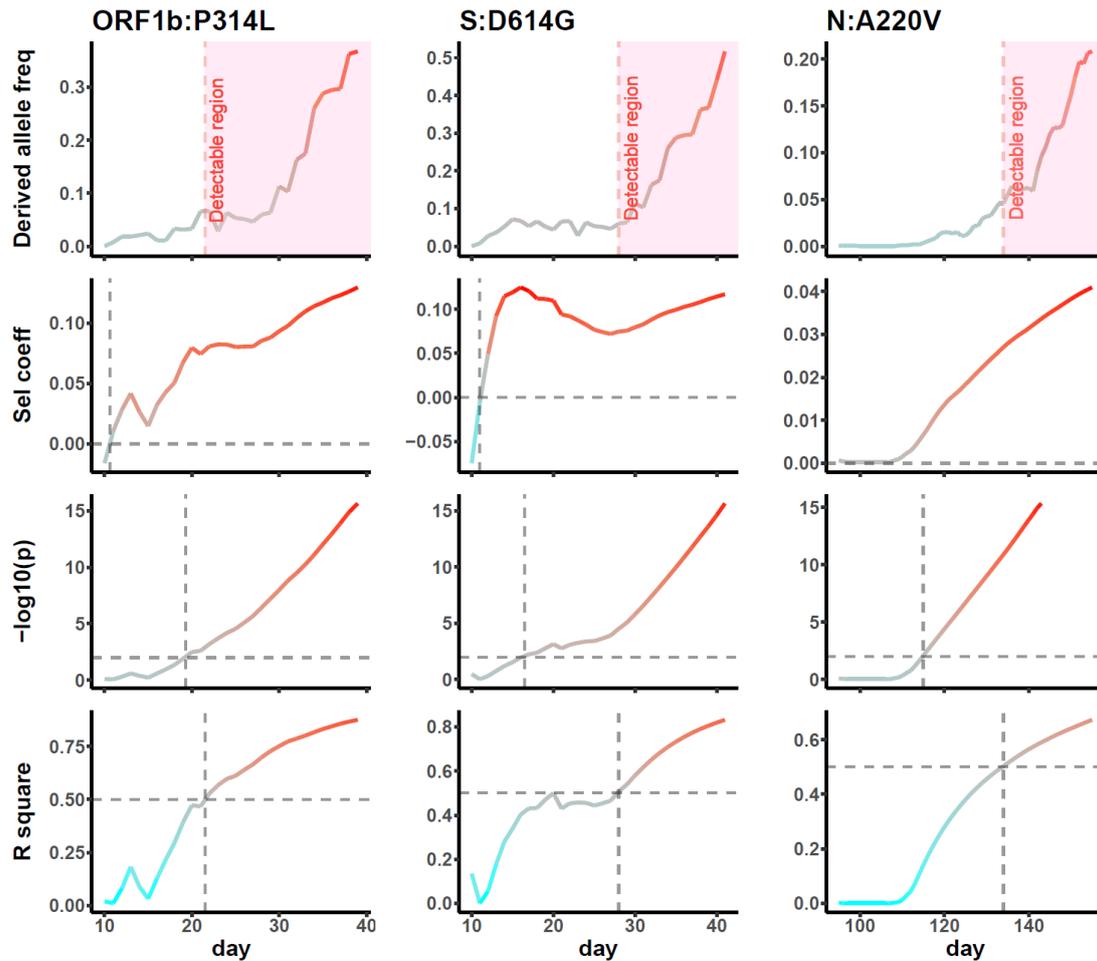
194 Portuguese: Huang Li<sup>7</sup>, Jingxuan Yan<sup>7</sup>, Chenglin Zhu<sup>7</sup>, Wenjing Liu<sup>7</sup>

- 195 Russian: Yuhong Guan<sup>7</sup>, Huanxiao Song<sup>7</sup>  
196 Spanish: Qin Zhou<sup>7</sup>, Han Gao<sup>7</sup>, Jinglan He<sup>7</sup>, Tiantian Li<sup>7</sup>, Ruiwen Fei<sup>7</sup>, Shumei  
197 Zhang<sup>7</sup>  
198 French: Yuyuan Guo<sup>7</sup>



199  
200  
201  
202  
203  
204  
205

**Figure 1. Timely-update and visualization framework of the Coronavirus GenBrowser.** The pre-analyzed genomic variant data can be freely accessed via <https://bigd.big.ac.cn/ncov/apis/>.



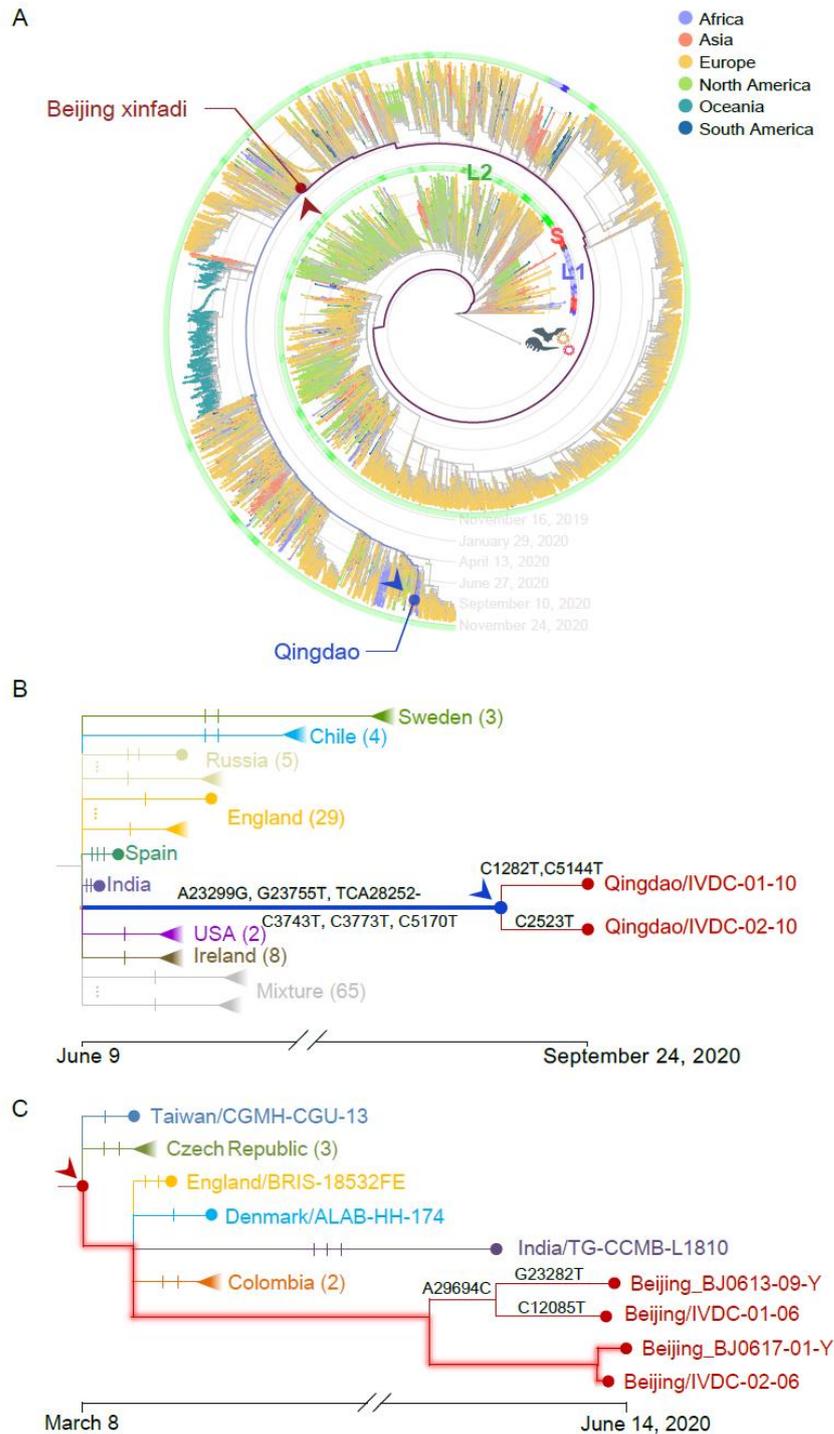
206

207

208 **Figure 2. Putative advantageous variants of SARS-CoV-2.** The  $x$ -axis displays  
209 number of days since the first appearance of derived allele in the global viral  
210 population. Predicted adaptation is marked in pink. Dashed gray crossings denote  
211 meaningful top right corners with a positive selection coefficient,  $p < 0.01$ , and  
212  $R^2 > 50\%$ .

213

214



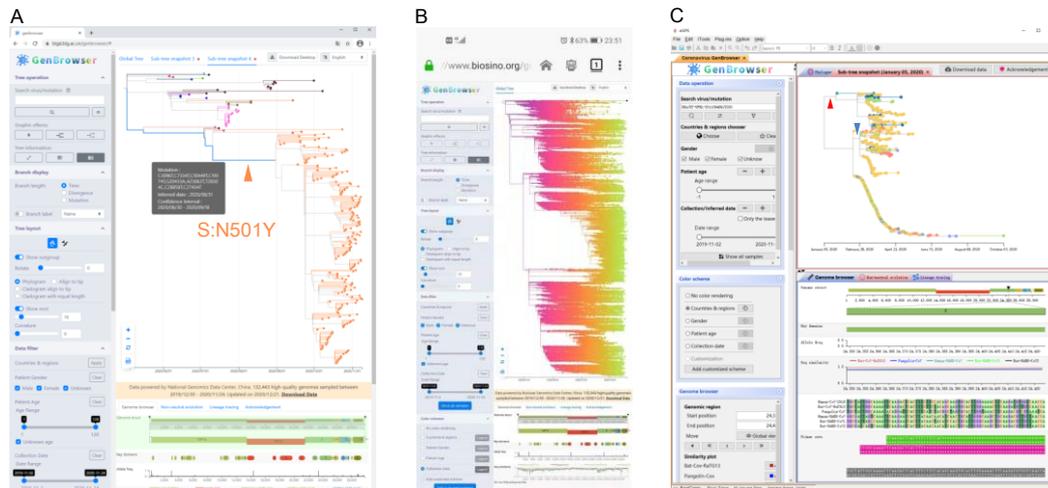
215  
216  
217 **Figure 3. Global and zoomed views of lineages associated with Qingdao and**  
218 **Beijing outbreaks.**

- 219 A) The lineages of traced targets are shown in blue and dark-red lines. The tree of  
220 132,443 viral strains was used. L/S lineage types<sup>30</sup> are marked with an outside  
221 circle.
- 222 B) Qingdao/IVDC-01-10 and Qingdao/IVDC-02-10 were the two SARS-CoV-2  
223 strains collected on September 24, 2020 from two dock workers in Qingdao,  
224 China. The query strain (env/Qingdao/IVDC-011-10) was found on an outer

225 packaging of cold-chain products on October 7, 2020 The environmental strain,  
226 marked with a blue solid circle with an arrow head, was found to be identical to  
227 the most recent common ancestor of the two strains from the two dock workers.  
228 Each notch of the branches represents a mutation. Mutations of the Qingdao  
229 strains are indicated.

230 C) The ancestral viral strain found in early March 2020 is marked with a dark-red  
231 solid circle with an arrow head. This strain is identical to the two strains  
232 (Beijing/IVDC-02-06, Beijing/BJ0617-01-Y) collected from two Xinfadi cases on  
233 June 11 and 14, 2020. The branches with no mutations are highlighted.

234



235

236

237 **Figure 4. Detection of non-neutral evolution of SARS-CoV-2 and tree**  
238 **visualization with CGB.**

239 A) Web-based CGB tree visualization of an accelerated lineage in the UK, out of  
240 132,443 SARS-CoV-2 genomic sequences, with the desktop version of Google  
241 Chrome.

242 B) Web-based CGB tree visualization of 132,443 genomes with the Android version  
243 of Firefox.

244 C) Tree visualization of a lineage (USA/UT-UPHL-201109489/2020) with the mostly  
245 reduced evolutionary rate and its neighbors with desktop standalone CGB. There  
246 are only two mutations (A20268G, red arrow head; C15324T, blue arrow head)  
247 happened in 966 strains within nearly 9 months.

248

249

## 250 References

251

- 252 1 Zhu, N. *et al.* A novel coronavirus from patients with pneumonia in China, 2019. *N Engl J*  
253 *Med* **382**, 727-733 (2020).
- 254 2 Lu, R. *et al.* Genomic characterisation and epidemiology of 2019 novel coronavirus:  
255 implications for virus origins and receptor binding. *Lancet* **395**, 565-574 (2020).
- 256 3 Wu, F. *et al.* A new coronavirus associated with human respiratory disease in China. *Nature*  
257 (2020).
- 258 4 Boni, M. F. *et al.* Evolutionary origins of the SARS-CoV-2 sarbecovirus lineage responsible  
259 for the COVID-19 pandemic. *Nature Microbiology* (2020).
- 260 5 Lan, J. *et al.* Structure of the SARS-CoV-2 spike receptor-binding domain bound to the ACE2  
261 receptor. *Nature* **581**, 215-220 (2020).
- 262 6 Hao, X. *et al.* Reconstruction of the full transmission dynamics of COVID-19 in Wuhan.  
263 *Nature* (2020).
- 264 7 Kissler, S. M., Tedijanto, C., Goldstein, E., Grad, Y. H. & Lipsitch, M. Projecting the  
265 transmission dynamics of SARS-CoV-2 through the postpandemic period. *Science* **368**,  
266 860-868 (2020).
- 267 8 Scudellari, M. The pandemic's future. *Nature* **584**, 22-25 (2020).
- 268 9 Hadfield, J. *et al.* Nextstrain: real-time tracking of pathogen evolution. *Bioinformatics* **34**,  
269 4121-4123 (2018).
- 270 10 Zhao, W.-M. *et al.* The 2019 novel coronavirus resource. *Yi Chuan* **42**, 212-221 (2020).
- 271 11 Shu, Y. L. & McCauley, J. GISAID: Global initiative on sharing all influenza data - from  
272 vision to reality. *Eurosurveillance* **22**, 2-4 (2017).
- 273 12 Elbe, S. & Buckland-Merrett, G. Data, disease and diplomacy: GISAID's innovative  
274 contribution to global health. *Glob Chall* **1**, 33-46 (2017).
- 275 13 Zhang, Z. *et al.* Database resources of the National Genomics Data Center in 2020. *Nucleic*  
276 *Acids Res* **48**, D24-D33 (2020).
- 277 14 Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7:  
278 Improvements in performance and usability. *Mol Biol Evol* **30**, 772-780 (2013).
- 279 15 Zhou, P. *et al.* A pneumonia outbreak associated with a new coronavirus of probable bat origin.  
280 *Nature* (2020).
- 281 16 Lam, T. T.-Y. *et al.* Identifying SARS-CoV-2-related coronaviruses in Malayan pangolins.  
282 *Nature* **583**, 282-285 (2020).
- 283 17 Hartigan, J. A. Minimum mutation fits to a given tree. *Biometrics* **29**, 53-65 (1973).
- 284 18 Zhao, Z. M. *et al.* Moderate mutation rate in the SARS coronavirus genome and its  
285 implications. *BMC Evol Biol* **4**, 21 (2004).
- 286 19 Cotten, M. *et al.* Spread, circulation, and evolution of the Middle East respiratory syndrome  
287 coronavirus. *Mbio* **5** (2014).
- 288 20 Li, X. *et al.* Evolutionary history, potential intermediate animal host, and cross-species  
289 analyses of SARS-CoV-2. *J Med Virol* (2020).
- 290 21 Flynn, J. A. *et al.* Exploring the coronavirus pandemic with the WashU Virus Genome  
291 Browser. *Nat Genet* **52**, 986-1001 (2020).
- 292 22 Fernandes, J. D. *et al.* The UCSC SARS-CoV-2 Genome Browser. *Nat Genet* **52**, 986-991  
293 (2020).

- 294 23 Korber, B. *et al.* Tracking changes in SARS-CoV-2 Spike: Evidence that D614G increases  
295 infectivity of the COVID-19 virus. *Cell* **182**, 812-827 (2020).
- 296 24 Hodcroft, E. B. *et al.* Emergence and spread of a SARS-CoV-2 variant through Europe in the  
297 summer of 2020. *medRxiv*, doi:<https://doi.org/10.1101/2020.10.25.20219063> (2020).
- 298 25 Xing, Y., Wong, G. W. K., Ni, W., Hu, X. & Xing, Q. Rapid response to an outbreak in  
299 Qingdao, China. *N Engl J Med*, doi:10.1056/NEJMc2032361 (2020).
- 300 26 Zhang, Y. *et al.* Genomic characterization of SARS-CoV-2 identified in a reemerging  
301 COVID-19 outbreak in Beijing's Xinfadi market in 2020. *Biosaf Health*,  
302 doi:<http://dx.doi.org/10.1016/j.bsheal.2020.08.006> (2020).
- 303 27 Pang, X. *et al.* Cold-chain food contamination as the possible origin of Covid-19 resurgence in  
304 Beijing. *Natl Sci Rev*, doi:10.1093/nsr/nwaa264 (2020).
- 305 28 Bedford, T. *et al.* Cryptic transmission of SARS-CoV-2 in Washington state. *Science* **370**,  
306 571-575 (2020).
- 307 29 Yu, D. *et al.* eGPS 1.0: comprehensive software for multi-omic and evolutionary analyses.  
308 *Natl Sci Rev* **6**, 867-869 (2019).
- 309 30 Tang, X. *et al.* On the origin and continuing evolution of SARS-CoV-2. *Natl Sci Rev* (2020).
- 310