

26 **Abstract**

27 Genomic epidemiology is important to study the COVID-19 pandemic and more than
28 two million SARS-CoV-2 genomic sequences were deposited into public databases.
29 However, the exponential increase of sequences invokes unprecedented bioinformatic
30 challenges. Here, we present the Coronavirus GenBrowser (CGB) based on a highly
31 efficient analysis framework and a movie maker strategy. In total, 1,002,739 high
32 quality genomic sequences with the transmission-related metadata were analyzed and
33 visualized. The size of the core data file is only 12.20 MB, efficient for clean data
34 sharing. Quick visualization modules and rich interactive operations are provided to
35 explore the annotated SARS-CoV-2 evolutionary tree. CGB binary nomenclature is
36 proposed to name each internal lineage. The pre-analyzed data can be filtered out
37 according to the user-defined criteria to explore the transmission of SARS-CoV-2.
38 Different evolutionary analyses can also be easily performed, such as the detection of
39 accelerated evolution and on-going positive selection. Moreover, the 75 genomic
40 spots conserved in SARS-CoV-2 but non-conserved in other coronaviruses were
41 identified, which may indicate the functional elements specifically important for
42 SARS-CoV-2. The CGB not only enables users who have no programming skills to
43 analyze millions of genomic sequences, but also offers a panoramic vision of the
44 transmission and evolution of SARS-CoV-2.

45

46 **Keywords**

47 Coronavirus GenBrowser; SARS-CoV-2; Genomic epidemiology; Transmission;
48 Evolution

49

50

51 **Introduction**

52 Real-time tracking of the transmission and evolution of the severe acute respiratory
53 syndrome coronavirus 2 (SARS-CoV-2) is essential for public health during the
54 COVID-19 pandemic [1]. Since January 2020 more than two million genomic
55 sequences have been deposited into public databases, such as National Center for
56 Biotechnology Information (NCBI) GenBank [2], Global Initiative on Sharing All
57 Influenza Data (GISAID) [3, 4]. The exponential increase of genomic sequences
58 provides a great opportunity to monitor the transmission and evolution of
59 SARS-CoV-2 but invokes unprecedented bioinformatic challenges.

60
61 Several web browsers have been developed to analyze the genomic data and track the
62 COVID-19 pandemic. The UCSC SARS-CoV-2 Genome Browser was derived from
63 the well-established UCSC genome-browser for visualization of nucleotide and protein
64 sequences, sequence conservations, and many other properties of wild-type and
65 variants of SARS-CoV-2 [5]. The WashU Virus Genome Browser provides
66 Nextstrain-based phylogenetic-tree view and genomic-coordinate, track-based view of
67 genomic features of viruses [6]. The pathogen genomics platform Nextstrain allows
68 analysis of genomic sequences of approximately 4,000 strains of SARS-CoV-2 and
69 investigation of its evolution [7], which cannot timely analyze millions of increasing
70 genomic sequences. Therefore, new approaches are essential and indispensable to
71 enable users easier to explore the large amount of SARS-CoV-2 genomic sequences.

72
73 In this study, we developed the Coronavirus GenBrowser (CGB). All the high-quality
74 genomic sequences and the associated transmission-related metadata were timely
75 analyzed to provide the latest panoramic view of the pandemic. To investigate a local
76 transmission, the data can be easily filtered according to countries, regions, keywords,
77 and the collection date of viral strains. Thus, even if users have no any programming
78 skills, the CGB enables them to efficiently explore millions of SARS-CoV-2 genomic

79 sequences and monitor the global/local transmission and evolution of SARS-CoV-2.
80 All the pre-analyzed genomic mutations and the associated metadata can be easily
81 downloaded, re-analyzed and re-shared. Since a cleaned genome alignment with
82 almost no ambiguous nucleotide sites can be easily re-constructed from the CGB core
83 data file, the CGB may provide a great convenience for the society to study viral
84 evolution further.

85

86 **Material and methods**

87 Brief descriptions of material and methods are included in this section. Detailed
88 descriptions are provided as Supplemental Materials.

89

90 **Data quality control and distributed genome alignments**

91 SARS-CoV-2 genomic variations were obtained from the 2019nCoV database [8]
92 established by China National Center for Bioinformation (CNCB) [9], as an integrated
93 resource based on Global Initiative on Sharing All Influenza Data (GISAID) [3, 4],
94 National Center for Biotechnology Information (NCBI) GenBank [2], China National
95 GeneBank DataBase (CNGBank) [10], the Genome Warehouse (GWH) [11], and the
96 National Microbiology Data Center (NMDC, <https://nmdc.cn/>). Detailed information
97 on this database is available at https://bigd.big.ac.cn/ncov/release_genome. All
98 SARS-CoV-2 strains were isolated from humans, and quality control was applied to
99 obtain high-quality SARS-CoV-2 genomic sequences (Figure S1, S2). Because of the
100 explosion in SARS-CoV-2 genomic data, the distributed alignment system was
101 developed to enable daily update (Figure 1), which reduces the total alignment time
102 complexity to $\mathcal{O}(n)$, where $\mathcal{O}(\cdot)$ is a linear function, and n is the number of viral
103 strains.

104

105 **Reconstruction and timely update of the annotated phylogenetic tree**

106 Before October 2020 all high-quality sequences in distributed alignments were
107 analyzed as a whole and used to reconstruct the evolutionary tree. After October 2020
108 the efficiency of this process became too low to perform timely updates. Therefore,
109 new trees were reconstructed by appending new sequences to existing tree.
110 Ambiguous and missing nucleotides were imputed by incorporation of the
111 neighboring lineages (Figure S3) and mutations in strains of each branch are
112 recapitulated according to the principle of parsimony [12, 13]. A highly effective
113 maximum-likelihood method (TreeTime) is used to determine the dates of internal
114 nodes [14] as it allows fast inference by “the post- and pre-order traversals” with
115 tabulated key values for back tracing. This algorithm was implemented in the CGB
116 with very minor revisions (Figures S4, and S5). The genome-wide mutation rate was
117 also timely calculated.

118

119 **Displaying SARS-CoV-2 genomic mutations in tree-based format**

120 Similar to NextStrain [7] and the WashU Virus Genome Browse [6], the CGB uses a
121 tree-based file format to store SARS-CoV-2 genomic mutations. The head of the core
122 data file contains the data version, the updated date, the genomic region analyzed, and
123 the mutation rate estimated for each gene. The core file also contains information on
124 collection date, gender and age of patient, location for each strain, mutations, and
125 inferred date for each internal node. To allow fast access to the data, redundant
126 information has been minimized.

127

128 **Visualization of the huge evolutionary tree by movie-maker strategy**

129 When visualizing the huge evolutionary tree, the most of lineages are invisible because
130 lineages overlap each other. Thus, a movie-maker strategy was implemented to
131 efficiently visualize huge evolutionary tree, in which only surface lines will be painted.
132 If the tree is zoomed in, only the visible sub-area of the tree will be painted. Using this
133 strategy, millions of lineages can be visualized effectively (Figure S6).

134

135 **CGB binary nomenclature and data searching**

136 A number of different naming systems have been proposed [15, 16], but these systems
137 only name a few internal nodes or branches. As there are a large number of internal
138 nodes on the huge evolutionary tree, the CGB binary nomenclature was developed
139 following the most recent common ancestor (MRCA) concept (Figure S11) to obtain
140 CGB ID for each node. CGB ID can be used to search a specific lineage. Isolate
141 names, accession numbers, and mutations are also searchable.

142

143 **Mutation analysis**

144 A root-to-tip linear regression method [17] was used to estimate the mutation rate of
145 SARS-CoV-2. For each strain with a different collection date in a tip-dated time tree,
146 the number of mutations, including that of recurrent mutations, was counted
147 subsequent to the appearance of MRCA. To avoid the effect of recombination,
148 recombination flag is labeled for each mutation by analyzing hybrid genomic
149 structure (Figure S12).

150

151 **Lineage tracing**

152 For lineage tracing, genomic sequences of SARS-CoV-2 strains collected from patients
153 or environments are used as the queries. These query sequences should be aligned
154 with the reference genomic sequence of SARS-CoV-2 (GenBank accession number:
155 NC_045512) [18] A very fast algorithm was implemented to count the difference
156 between a query sequence and the genomic sequence of a node. For one query, nodes
157 with the least difference are considered as its candidate targets.

158

159 **Detection of branch-specific accelerated evolution of SARS-CoV-2**

160 To detect branch-specific accelerated evolution, each internal branch of the
161 SARS-CoV-2 tree was examined. For each internal branch, the observed number of

162 mutations of the i -th gene ($\gamma_{obs,i}$) was compared with the expected number of mutations
163 of the same gene ($\gamma_{exp,i}$). The significance level of acceleration was determined by
164 Poisson probability [19, 20]. It is a one-tailed test. The condition $t > 10$ (days) was
165 used for detection of branch-specific accelerated evolution of SARS-CoV-2.

166

167 **Detection of on-going selection of SARS-CoV-2**

168 To detect on-going positive selection, allele frequency trajectory with an S-shaped
169 curve was examined (Table S3, Figure S13, S14). To reduce the impact of hitchhiking
170 by neutral mutation, only non-synonymous mutations were analyzed, although
171 non-coding mutations [21] can also be beneficial.

172

173 **Local analysis for new genomic data of SARS-CoV-2**

174 The public global genomic data can be freely downloaded and analyzed together with
175 new genomic data before these new sequences are integrated into the CGB by the
176 CGB team members. This function ensures that a timely analysis can be easily
177 performed.

178

179 **Data source**

180 For genomic sequence alignments, high quality SARS-CoV-2 genomic variations were
181 obtained from the 2019nCoV database [8, 22], which is an integrated resource based
182 on GenBank, GISAID [3, 4], China National GeneBank DataBase (CNGBdb) [10], the
183 Genome Warehouse (GWH) [11], and the National Microbiology Data Center (NMDC,
184 <https://nmdc.cn/>).

185

186 **Results**

187 **The construction of a million level evolutionary tree**

188 After quality control, 1,002,739 high-quality genomic sequences were obtained for
189 subsequent analyses. The number of identified high- and low-quality genomes in each

190 month is shown (Figure S2). To allow timely analysis of a large number of sequences,
191 we first solved the problem that all viral genomic sequences have to be re-aligned when
192 nucleotide sequences of new genomes become available. This is extremely
193 time-consuming. With the distributed alignment system (Figure 1), we dramatically
194 reduced the total time required for the alignment. We also built the evolutionary tree on
195 the existing tree with new genomic data in order to reduce the complexity of tree
196 construction. With these modifications, a tremendous evolutionary tree can be
197 reconstructed for each update, millions of SARS-CoV-2 genomic sequences can be
198 timely analyzed, and data can be easily shared, reanalyzed and re-constructed (Figure
199 1).

200

201 For the huge evolutionary tree, mutations on each branch were identified according to
202 the principle of parsimony [12, 13] and the dates of internal nodes were inferred with
203 very minor revisions of a highly effective maximum-likelihood method (TreeTime)
204 [14]. The pre-analyzed genomic mutations of SARS-CoV-2 and the associated
205 metadata are shared to the general public in a tree-based CGB format. The size of
206 distributed alignments is 30.28 GB for the 1,002,739 SARS-CoV-2 genomic sequences.
207 The tree-based data format allows the compression ratio to reach 2,541:1, meaning that
208 the size of compressed core file containing the pre-analyzed genomic mutations and
209 associated metadata is as small as 12.20 MB with zip compression (Figure 1).

210 Whenever necessary, cleaned genomic sequences with almost no ambiguous
211 nucleotide sites can be reconstructed for viral isolates. Thus, this approach ensures
212 low-latency access to the data and enables fast sharing and re-analysis of a large
213 number of SARS-CoV-2 genomic variants.

214

215 **Highly efficient visualization of the tree and tracks**

216 To efficiently visualize the results, a movie-maker strategy was implemented for
217 painting the evolutionary tree, only elements shown on the screen and visible to the user

218 are painted. This design makes the visualization process highly efficient, and the
219 evolutionary tree of more than one million strains can be visualized. It takes about one
220 second for the visualization process in different operation systems (Table S1).

221

222 **The convenience of tree operation**

223 The CGB is also a highly efficient platform to search or filter variants based on
224 transmission-related metadata (Figure 2). Useful interactive functionalities were
225 developed to navigate users through the huge tip-dated evolutionary tree. First, users
226 can easily search internal branches or variants with certain mutations, or isolate names
227 of virus. There are 400,298 internal branches in the evolutionary tree ($n =$
228 1,002,739), and each branch has been named by CGB binary nomenclature (*i.e.*, CGB
229 ID) (Figure S11) and is searchable. Thus, different Variants of Concerns (VOCs) can
230 be easily identified and visualized on the huge evolutionary tree (Table 1). Second,
231 users can easily filter out the data according to the collection date, the country/region,
232 and the gender and age groups. Third, the visualization of a sub-clade in another tab is
233 allowed. Forth, different annotations are provided to mark the clades of interest. At last,
234 coordinated annotation tracks are provided to show genome structure and key domains,
235 allele frequencies, sequence similarity between various coronavirus, multi-genome
236 alignment and primer sets for detection of SARS-CoV-2 (Figure S7-10). All those
237 features have brought great convenience for users.

238

239 **Transmission case studies through CGB**

240 Based on the large data volume and user-friendly CGB, many analyses can be quickly
241 conducted. Take the spread of VOCs in India as an example. Among the 1,002,739
242 strains, there were 3,349 ones sampled in India. Nearly all major SARS-CoV-2
243 lineages can be found in India in different stages of the pandemic (Figure 3). By
244 searching the CGB IDs of the VOCs (Table 1), the clades of VOCs were identified. In
245 total, there were 464 ($464/3,349=13.85\%$) Delta strains, 185 ($185/3,349=5.52\%$)

246 Alpha strains, and 11 (11/3,349=0.32%) Beta strains in the Indian sample. The ratios
247 change when only considering recent viral strains collected after Apr 01, 2021. There
248 were 589 India samples after date filtering. Among them, there were 364
249 (364/589=61.80%) Delta strains, 55 (55/589=9.34%) Alpha (B.1.1.7) strains and one
250 (1/589=0.17%) Beta strains. Thus, the Delta variant increased more rapidly than
251 others, and the most recent infections in India are caused by the Delta (B.1.617.2)
252 VOC.

253

254 Another three examples were provided to show that the CGB is an efficient platform to
255 investigate local and global transmission of COVID-19 (Figure 4). To trace the origin
256 of a local COVID-19 outbreak, the lineage tracing was implemented in the CGB. The
257 closest nodes were revealed for the three outbreaks in China during this year which
258 indicates different origin of the three outbreaks. Their neighboring strains can be
259 viewed individually and further investigated. The analysis is extremely fast and can
260 be performed on a desktop computer. Therefore, the CGB is a highly efficient
261 platform to investigate the origin of a local COVID-19 outbreak.

262

263 **Mutation analysis**

264 The CGB also estimates the mutation rate of whole genome and each gene (Table S2).
265 Applying 1,002,739 genomic sequences, the estimated genome-wide mutation rate
266 is 1.0794×10^{-3} per nucleotide per year. The mutation rate is variable for different
267 genes (Table S2).

268

269 We also found that the mutation rate could be different among sites. Using the CGB
270 core data file, we conducted a 10-base sliding window analysis with a sliding step of
271 one base and identified fine-scaled mutation cold spots along the viral genome,
272 indicating the genomic regions with mutation rate significantly lower than the average
273 mutation rate of the entire genome. In total, 657,074 (recurrent) mutations were

274 identified and 868 mutation cold spots were found with a false discovery rate (FDR)
275 corrected P -value < 0.01 (Figure 5, Supplemental excel file). The coldest spot is
276 located in ORF1ab, which encodes nsp13 helicase (nucleotides 16, 294 – 16,307)
277 (FDR corrected P -value = 4.79×10^{-46}). Interestingly, it has been found that
278 sequence conservation is restricted to ORF1ab:nsp10-13 among 14 coronaviruses [23,
279 24]. It indicates that nsp13 helicase might be essential for coronaviruses and
280 SARS-CoV-2. Moreover, among the 868 mutation cold spots, there are 75 conserved
281 spots in SARS-CoV-2, but not conserved among other coronaviruses. These
282 SARS-CoV-2 specific conserved elements may play key roles for SARS-CoV-2
283 specific functions.

284

285 **The detection of accelerated evolution**

286 The CGB provides a module to detect branch-specific accelerated evolution of
287 SARS-CoV-2. We found that within 186,746 internal branches with $t > 10$ (days) on
288 the evolutionary tree ($n = 1,002,739$), 70 branches were detected to have a
289 genome-wide accelerated evolution (FDR corrected $P < 0.05$), 332 branches were
290 detected to have an accelerated evolution of ORF1ab (FDR corrected $P < 0.05$), and
291 two branches was found to have an accelerated evolution of the spike gene (FDR
292 corrected $P < 0.05$) (Supplemental excel file). These evolution-accelerated variants can
293 be used for future studies.

294

295 **The detection of on-going positive selection**

296 The mutation frequency trajectory for each mutation can be easily visualized by using
297 the CGB. It has a module to detect on-going positive selection based on S-shaped
298 frequency trajectory of a selected allele (Figures S13, S14). It has been shown that the
299 SARS-CoV-2 variant with G614 spike protein has a fitness advantage [25]. Our
300 analysis using the CGB confirmed this finding even when the G614 frequency was
301 very low ($< 10\%$) (Figure S15), indicating that the CGB can detect putative

302 advantageous variants before they become widely spread. The CGB also predicted an
303 increase in the frequency of S:p.P681R of the Delta VOC (Figure S15), suggesting
304 that variants with the mutation may be advantageous. S:p.P681R is located on the
305 spike S1/S2 cleavage site, and another mutation (S:p.P681H) on the same position has
306 been found to be advantageous [26] and may interact with other mutations [21] in the
307 Alpha VOC. Based on 1,002,739 samples, the CGB detected 13 putative
308 advantageous mutations in the spike protein (Table S4). However, as an increase in
309 mutation frequency could be due to sampling bias and epidemiological factors [25],
310 putative advantageous variants should be closely monitored.

311

312 **Conclusion**

313 In this study, we developed an effective surveillance tool for the transmission and
314 evolution of SARS-CoV-2. A highly efficient visualization module is established, and
315 rich interactive operations are allowed to explore the annotated evolutionary tree. We
316 investigated four local COVID-19 outbreaks by the searching and lineage tracing
317 functionalities. We also implemented a new method to detect on-going positive
318 selection for each viral nonsynonymous mutation. The branches with accelerated or
319 reduced evolutionary rate are identified to provide a real time tracking on the change
320 of evolutionary rate, which could reveal epidemic factors affecting the viral
321 transmission.

322

323 The CGB provides an efficient way for clean data sharing. It could be difficult for
324 researchers to download all the raw genomic sequences, perform data quality control
325 and analyze the large amount of data by their own. By examining literatures related to
326 this topic [27, 28], most of those studies have similar methods to prepare alignments,
327 build phylogenetic tree, infer mutations and date each internal node. The data
328 preparation procedure is tedious but may require some skills and is often time
329 consuming when the sample size is extremely large. Nevertheless, by downloading

330 the pre-analyzed CGB data files, users can have all those clean data in a few minutes.
331 To promote a timely analysis of newly sequenced genomic data, users can perform a
332 local analysis to analyze these new data, together with the public genomic data
333 globally sampled. Thus, the CGB provides a convenient way to study the evolution of
334 SARS-CoV-2 and monitor its transmission.

335

336 During the analysis, we noticed that a very small percentage of sequences have
337 abnormal collection dates which could severely skew the evolutionary tree. After
338 examining all possible reasons, it is likely due to that the year of collection date was
339 incorrectly filled in the most cases. We then deleted these sequences although their
340 sequence quality is high. Thus, we would suggest that researchers could pay attention
341 on the year of collection date when submitting their sequences.

342

343 The public science education is extremely important for the anti-epidemic to show
344 that SARS-CoV-2 has been evolving. Therefore, a web-based CGB was also
345 developed. It is a simplified version of CGB that provides a convenient way to access
346 the data via a web browser, such as Google Chrome, Firefox, and Safari (Figure S6).
347 For educational purpose, nine language versions (Chinese, English, German, Japanese,
348 French, Italian, Portuguese, Russian, and Spanish) are available. The web-based CGB
349 package can be downloaded and reinstalled on any websites. Two pre-installed
350 websites are provided (<https://www.biosino.org/genbrowser/> and
351 <https://ngdc.cncb.ac.cn/genbrowser/>). For the scientific community, it is highly
352 recommended to download the eGPS software to use the CGB
353 (http://www.egps-software.net/egpscloud/eGPS_Desktop.html). Overall, the CGB is
354 frequently updated which provides a timely panoramic vision of the global and local
355 transmission and evolution of SARS-CoV-2.

356

357

358 **Funding**

359 This work was supported by a grant from the National Key Research and Development
360 Project [No. 2020YFC0847000]. Funding for open access charge: National Key
361 Research and Development.

362

363 **Acknowledgements**

364 We thank Ya-Ping Zhang for providing valuable advice and encouragement and the
365 researchers who generated and deposited sequence data of SARS-CoV-2 in GISAID,
366 GenBank, CNGBdb, GWH, and NMDC making this study possible.

367

368 **Data availability**

369 All timely-updated data are freely available at <https://bigd.big.ac.cn/ncov/apis/>. The
370 desktop standalone version (Figure S6A) provides the full function of CGB and has a
371 plug-in module for the eGPS software
372 (http://www.egps-software.net/egpscloud/eGPS_Desktop.html) [29].

373

374 **Members of the language translation team**

375 German: Ning He⁶, Jing Lv⁶, Ting Peng⁶

376 Italian: Ting Zhou⁶, Nan Yang⁶, Siyi Hou⁶

377 Portuguese: Huang Li⁶, Jingxuan Yan⁶, Chenglin Zhu⁶, Wenjing Liu⁶

378 Russian: Yuhong Guan⁶, Huanxiao Song⁶

379 Spanish: Qin Zhou⁶, Han Gao⁶, Jinglan He⁶, Tiantian Li⁶, Ruiwen Fei⁶, Shumei Zhang⁶

380 French: Yuyuan Guo⁶

381

382 **Conflict of interest**

383 The authors declare no competing interests.

384

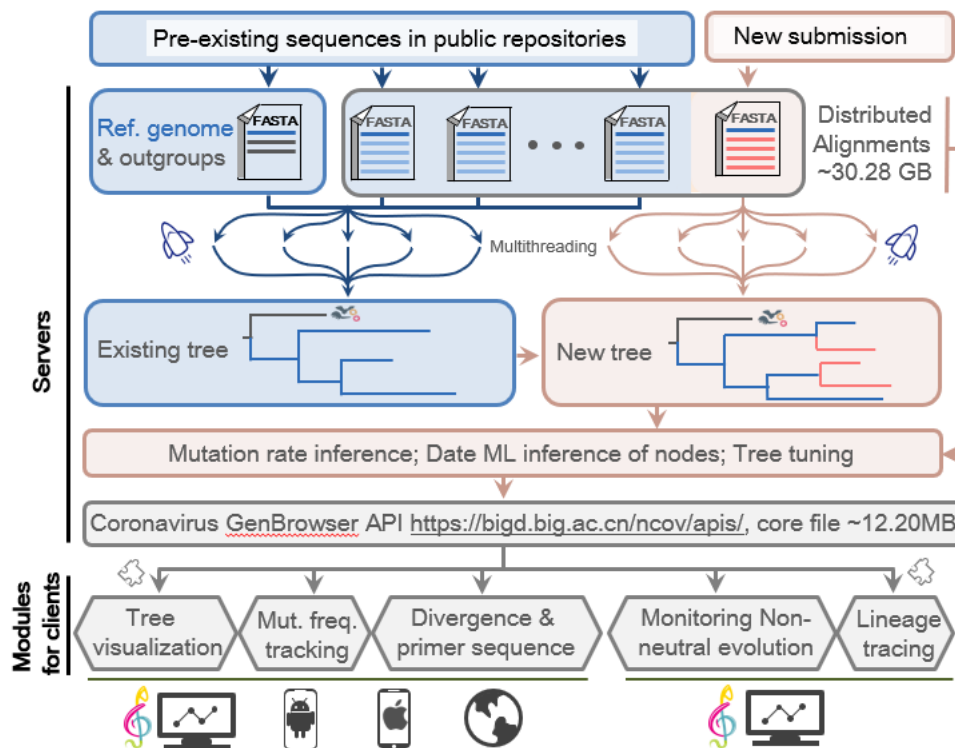
385

386

387 **References**

- 388 1. Fineberg HV, Wilson ME. Epidemic science in real time, *Science* 2009;324:987.
- 389 2. Sayers EW, Beck J, Bolton EE et al. Database resources of the National Center for
390 Biotechnology Information, *Nucleic Acids Res* 2021;49:D10-D17.
- 391 3. Shu YL, McCauley J. GISAID: Global initiative on sharing all influenza data -
392 from vision to reality, *Eurosurveillance* 2017;22:2-4.
- 393 4. Elbe S, Buckland-Merrett G. Data, disease and diplomacy: GISAID's innovative
394 contribution to global health, *Glob Chall* 2017;1:33-46.
- 395 5. Fernandes JD, Hinrichs AS, Clawson H et al. The UCSC SARS-CoV-2 Genome
396 Browser, *Nat Genet* 2020;52:986-991.
- 397 6. Flynn JA, Purushotham D, Choudhary MNK et al. Exploring the coronavirus
398 pandemic with the WashU Virus Genome Browser, *Nat Genet* 2020;52:986-1001.
- 399 7. Hadfield J, Megill C, Bell SM et al. Nextstrain: real-time tracking of pathogen
400 evolution, *Bioinformatics* 2018;34:4121-4123.
- 401 8. Zhao W-M, Song S-H, Chen M-L et al. The 2019 novel coronavirus resource,
402 *Hereditas (Beijing)* 2020;42:212-221.
- 403 9. Xue YB, Bao YM, Zhang Z et al. Database resources of the National Genomics
404 Data Center, China National Center for Bioinformation in 2021, *Nucleic Acids Res*
405 2021;49:D18-D28.
- 406 10. Chen F, You L, Yang F et al. CNGBdb: China National GeneBank DataBase,
407 *Hereditas (Beijing)* 2020;42:799-809.
- 408 11. Chen M, Ma Y, Wu S et al. Genome Warehouse: A public repository housing
409 genome-scale data, *Genomics Proteomics Bioinformatics* 2021.
- 410 12. Sankoff D. Minimal mutation trees of sequences., *SIAM J Appl Math*
411 1975;28:35-42.
- 412 13. Hartigan JA. Minimum mutation fits to a given tree, *Biometrics* 1973;29:53-65.
- 413 14. Sagulenko P, Puller V, Neher RA. TreeTime: Maximum-likelihood phylodynamic
414 analysis, *Virus Evol* 2018;4:vex042.
- 415 15. Tang X, Wu C, Li X et al. On the origin and continuing evolution of SARS-CoV-2,
416 *Natl Sci Rev* 2020;7:1012-1023.
- 417 16. Forster P, Forster L, Renfrew C et al. Phylogenetic network analysis of
418 SARS-CoV-2 genomes, *Proc Natl Acad Sci USA* 2020;117:9241-9243.
- 419 17. Bouckaert R, Vaughan TG, Barido-Sottani J et al. BEAST 2.5: An advanced
420 software platform for Bayesian evolutionary analysis, *PLoS Comput Biol* 2019;15.
- 421 18. Wu F, Zhao S, Yu B et al. A new coronavirus associated with human respiratory
422 disease in China, *Nature* 2020;579:265-269.
- 423 19. Ohta T, Kimura M. On the constancy of the evolutionary rate in cistrons, *J Mol*
424 *Evol* 1971;1:18-25.
- 425 20. Wang Y, Dai G, Gu Z et al. Accelerated evolution of an *Lhx2* enhancer shapes
426 mammalian social hierarchies, *Cell Res* 2020;30:408-420.
- 427 21. Yang J, Zhang G, Yu D et al. A Kozak-related non-coding deletion effectively
428 increases B.1.1.7 transmissibility, *bioRxiv* 2021.

- 429 22. Gong Z, Zhu J-W, Li C-P et al. An online coronavirus analysis platform from the
430 National Genomics Data Center, *Zool Res* 2020;41:705-708.
- 431 23. Ruan YJ, Wei CL, Ee LA et al. Comparative full-length genome sequence analysis
432 of 14 SARS coronavirus isolates and common mutations associated with putative
433 origins of infection, *Lancet* 2003;361:1779-1785.
- 434 24. He JF, Peng GW, Min J et al. Molecular evolution of the SARS coronavirus during
435 the course of the SARS epidemic in China, *Science* 2004;303:1666-1669.
- 436 25. Korber B, Fischer WM, Gnanakaran S et al. Tracking changes in SARS-CoV-2
437 Spike: Evidence that D614G increases infectivity of the COVID-19 virus, *Cell*
438 2020;182:812-827.
- 439 26. Rambaut A, Loman N, Pybus O et al. Preliminary genomic characterisation of an
440 emergent SARS-CoV-2 lineage in the UK defined by a novel set of spike
441 mutations, virological.org
442 2020:[https://virological.org/t/preliminary-genomic-characterisation-of-an-emerge](https://virological.org/t/preliminary-genomic-characterisation-of-an-emergent-sars-cov-2-lineage-in-the-uk-defined-by-a-novel-set-of-spike-mutations/563)
443 [nt-sars-cov-2-lineage-in-the-uk-defined-by-a-novel-set-of-spike-mutations/563](https://virological.org/t/preliminary-genomic-characterisation-of-an-emergent-sars-cov-2-lineage-in-the-uk-defined-by-a-novel-set-of-spike-mutations/563).
- 444 27. Hodcroft EB, Zuber M, Nadeau S et al. Spread of a SARS-CoV-2 variant through
445 Europe in the summer of 2020, *Nature* 2021.
- 446 28. Deng XD, Gu W, Federman S et al. Genomic surveillance reveals multiple
447 introductions of SARS-CoV-2 into Northern California, *Science*
448 2020;369:582-587.
- 449 29. Yu D, Dong L, Yan F et al. eGPS 1.0: comprehensive software for multi-omic and
450 evolutionary analyses, *Natl Sci Rev* 2019;6:867-869.
- 451 30. O'Toole Á, Hill V, Pybus OG et al. Tracking the international spread of
452 SARS-CoV-2 lineages B.1.1.7 and B.1.351/501Y-V2, *Wellcome Open Res*
453 2021;6:121.
- 454
- 455



456

457

458 **Figure 1. Timely updates of SARS-CoV-2 genomic data and visualization**

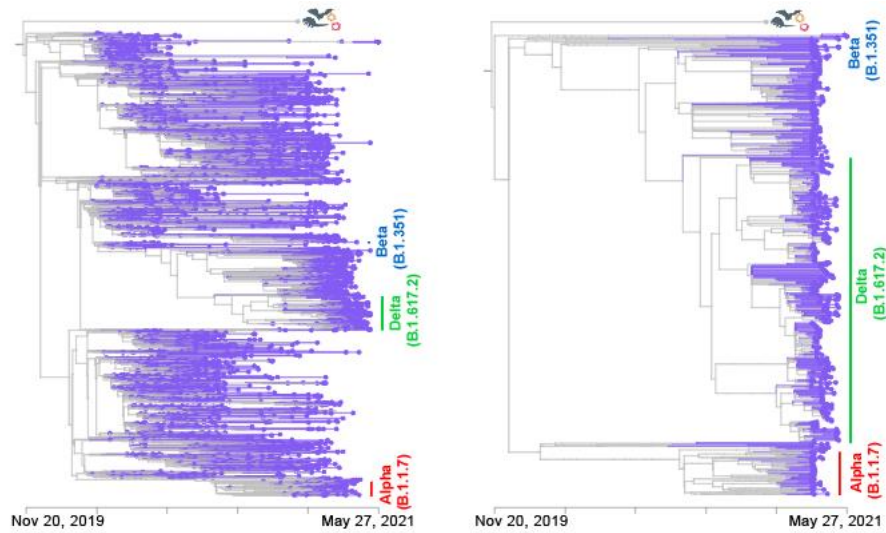
459 **framework of Coronavirus GenBrowser.**

460 The core file includes the pre-analyzed genomic mutations of SARS-CoV-2 and the

461 associated metadata. All timely-updated data can be freely accessed at

462 <https://bigd.big.ac.cn/ncov/apis/>.

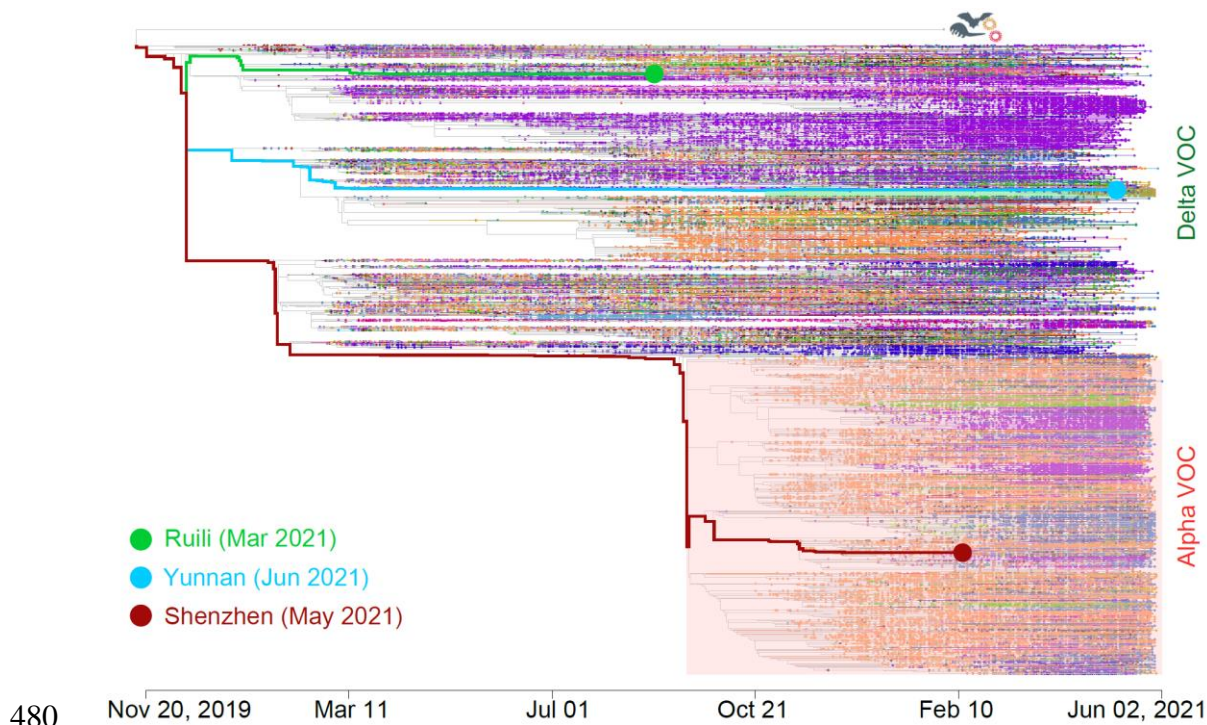
463



474

475 **Figure 3. Proportion change of VOCs in India.**

476 The tree on the left shows all India samples of 1,002,739 SARS-CoV-2 strains and the
477 tree on the right shows India samples after Apr 01, 2021. Different VOC are
478 annotated in different colors, Alpha (B.1.1.7) VOC in red, Delta (B.1.617.2) VOC in
479 green and Beta (B.1.351) VOC in blue, no Gamma are found in India.



480

481 **Figure 4. Tracing the origin of three local outbreak of COVID-19 in China.**

482 One sequence sampled during each outbreak is used as the query sequence which is
483 not included in the CGB dataset version used in analysis. Their closest targets were
484 marked with colored dot and the evolutionary paths were highlighted. The GISAID
485 IDs for the queries are EPI_ISL_1595852 (Ruili), EPI_ISL_2834004 (Yunnan), and
486 EPI_ISL_2405168 (Shenzhen).

487



488

489 **Figure 5. Manhattan plot of mutation cold spots in the genome of SARS-CoV-2.**

490 Results of genome-wide scan for mutation cold spots are shown in Manhattan plot of

491 significance against SARS-CoV-2 reference genomic locations. In total, 1,002,739

492 high quality genomic sequences were analyzed. Each dot represents one window.

493 *P*-values are FDR-corrected. The dotted red line denotes FDR-corrected *P*-value < 0.01.

494 Dots above the line represent mutation cold spots. Genomic structure and sequence

495 similarity between SARS-CoV-2 reference genome (NC_045512.2)[18] and the

496 genomes of five other coronaviruses are shown above the Manhattan plot.

497 **Table 1. CGB ID for Variant of Concern (VOC).**

WHO label	Pango lineage	Documented samples	CGB ID	Defining SNPs ^a
Alpha	B.1.1.7	United Kingdom, Sep-2020	CGB84017.91425	ORF1ab: T1001I, A1708D, I2230T, SGF3675- S: HV69-, Y144-, N501Y, A570D, P681H, T716I, S982A, D1118H ORF8: Q27*, R52I, Y73C N: D3L, S235F
Beta	B.1.351	South Africa, May-2020	CGB391494.393307	E: P71L N: T205I ORF1a: K1655N S: D80A, D215G, K417N, A701V, N501Y, E484K
Gamma	P.1	Brazil, Nov-2020	CGB222196.451180	ORF1ab: S1188L, K1795Q, SGF3675- S: L18F, T20N, P26S, D138Y, R190S, K417T, E484K, N501Y, H655Y, T1027I ORF3a: G174C ORF: E92K N: P80R
Delta	B.1.617.2	India, Oct-2020	CGB531065.580055	S: T19R, L452R, T478K, P681R, D950N ORF3a: S26L M: I82T ORF7a: V82A, T120I N: D63G, R203M, D377Y

498 ^a The information of defining SNPs were obtained from the Pango lineages[30].