

## An evolutionary portrait of the progenitor SARS-CoV-2 and its dominant offshoots in COVID-19 pandemic

Sudhir Kumar<sup>1,2,\*</sup>, Qiqing Tao<sup>1,2</sup>, Steven Weaver<sup>1,2</sup>, Maxwell Sanderford<sup>1,2</sup>, Marcos A. Caraballo-Ortiz<sup>1,2</sup>, Sudip Sharma<sup>1,2</sup>, Sergei L. K. Pond<sup>1,2,\*</sup>, and Sayaka Miura<sup>1,2,\*</sup>

<sup>1</sup>Institute for Genomics and Evolutionary Medicine, Temple University, Philadelphia, PA.

<sup>2</sup>Department of Biology, Temple University, Philadelphia, PA.

**\*Co-corresponding author:**

Sudhir Kumar ([s.kumar@temple.edu](mailto:s.kumar@temple.edu))

Sergei Pond ([spond@temple.edu](mailto:spond@temple.edu))

Sayaka Miura ([sayaka.miura@temple.edu](mailto:sayaka.miura@temple.edu))

## 1 Abstract

2 We report the likely most recent common ancestor of SARS-CoV-2 – the coronavirus that causes COVID-  
3 19. This progenitor SARS-CoV-2 genome was recovered through a novel application and advancement of  
4 computational methods initially developed to reconstruct the mutational history of tumor cells in a  
5 patient. The progenitor differs from the earliest coronaviruses sampled in China by three variants,  
6 implying that none of the earliest patients represent the index case or gave rise to all the human  
7 infections. However, multiple coronavirus infections in China and the USA harbored the progenitor  
8 genetic fingerprint in January 2020 and later, suggesting that the progenitor was spreading worldwide as  
9 soon as weeks after the first reported cases of COVID-19. Mutations of the progenitor and its offshoots  
10 have produced many dominant coronavirus strains, which have spread episodically over time.  
11 Fingerprinting based on common mutations reveals that the same coronavirus lineage has dominated  
12 North America for most of the pandemic. There have been multiple replacements of predominant  
13 coronavirus strains in Europe and Asia and the continued presence of multiple high-frequency strains in  
14 Asia and North America. We provide a continually updating dashboard of global evolution and  
15 spatiotemporal trends of SARS-CoV-2 spread (<http://sars2evo.datamonkey.org/>).

16

## 17 Main

18 Despite an unprecedented scope of global genome sequencing of Severe acute respiratory syndrome  
19 coronavirus 2 (SARS-CoV-2) and a multitude of phylogenetic analyses<sup>1-5</sup>, the early evolutionary history  
20 of SARS-CoV-2 remains unclear. Sophisticated investigations have found that traditional molecular  
21 phylogenetic analyses do not produce reliable evolutionary inferences about the early history of SARS-  
22 CoV-2 due to low sequence divergence, a limited number of phylogenetically informative sites, and the  
23 ubiquity of sequencing errors<sup>6-8</sup>. In particular, the root of the SARS-CoV-2 phylogeny remains elusive<sup>9,10</sup>  
24 because the closely-related non-human coronavirus (outgroups) more than 1,100 base differences  
25 from human SARS-CoV-2 genomes, as compared to fewer than 30 differences between human SARS-  
26 CoV-2 genomes' sequenced early on (December 2019 and January 2020)<sup>7,9-15</sup>. Without a reliable root  
27 of the SARS-CoV-2 phylogeny, one cannot accurately reconstruct the most recent ancestor sequence.  
28 Consequently, we cannot determine if any of the coronaviruses isolated to date carried the genome of  
29 the most recent common ancestor (progenitor) of all human SARS-CoV-2 infections. Knowing the  
30 progenitor genome will help us determine how close the earliest patients sampled in China represent  
31 are to "patient zero," i.e., the first case of human transmission.

32 The orientation and order of early mutations giving rise to common coronavirus variants will be misled  
33 if the earliest coronavirus isolates are incorrectly used to root the SARS-CoV-2 phylogenies<sup>3,16-18</sup>. The  
34 earliest investigations of COVID-19 patients and their coronaviruses' genomes already reported the  
35 presence of multiple variants<sup>19,20</sup>, and genomes of viral samples from December 2019 had as many as

36 five differences from each other. These observations require an explicit test of the assumption that one  
37 of the early sampled coronavirus genomes was the most recent common ancestor (progenitor) of all  
38 the strains infecting humans. Traditionally, the ancestral sequence of organisms is estimated by using  
39 a rooted phylogeny<sup>21,22</sup>. This ancestral sequence can then be compared with sequenced genomes to  
40 find the one that is most similar to that of the inferred progenitor and/or placed closest to the root in  
41 the phylogeny. However, as noted above, attempts using *ad hoc* and traditional methods are fraught  
42 with difficulties and have produced contradictory results<sup>9,10</sup>. Some methods also incorporate sampling  
43 times in phylogenetic inference, but they favor placing the earliest sampled genomes at or near the  
44 root of the tree<sup>10</sup>. This practice introduces a degree of circularity in testing the hypothesis that the  
45 earliest sampled genomes were ancestral because sampling time is used in the inference procedure.

## 46 Results and Discussion

### 47 A mutational order approach for SARS-CoV-2

48 We applied a mutation order approach (MOA) that directly reconstructs the ancestral sequence and  
49 the mutational history of genomes<sup>23–25</sup> without inferring a phylogeny as an intermediate step. MOA is  
50 often used to reconstruct the evolutionary history of tumor cells that evolve clonally and without  
51 recombination. This approach is well-suited for analyzing SARS-CoV-2 genomes because of their quasi-  
52 species evolutionary behavior (clonal) and because of the lack of evidence of significant recombination  
53 within human outbreaks, both of which preserve the collinearity of variants in genomes. This feature  
54 permits effective use of shared co-occurrence of variants in genomes, as well as the frequencies of  
55 individual variants, to infer mutational history, notwithstanding the presence of sequencing errors and  
56 other artifacts<sup>23,26</sup> (see *Methods*). We advanced MOA for application in the analysis of SARS-CoV-2  
57 genomes because the normal cell sequence in tumors provides a direct way to establish the ancestral  
58 (non-cancerous) genome. Such a direct ancestor is not available for coronaviruses in which the closest  
59 outgroup sequences are over 30-times more different than any two human strains. We also devised a  
60 bootstrap approach to place confidence limits on the inferred mutation order in which bootstrap  
61 replicate datasets are generated by sampling genomes with replacement (see *Methods*).

62 We analyzed two snapshots of the fast-growing collection of SARS-CoV-2 genomes to make inferences  
63 and assess the robustness of the inferred mutational histories to the growing genome collection,  
64 expanding at an unprecedented rate. The first snapshot was retrieved from GISAID<sup>27</sup> on July 7, 2020,  
65 and consisted of 60,332 genomes. Of these, 29,681 were selected because they were longer than the  
66 28,000 bases threshold imposed (29KG dataset) and did not include an excessive number of unresolved  
67 bases in any genomic regions. This second snapshot was acquired on October 12, 2020, from GISAID  
68 and contained 133,741 genomes, of which 68,057 genomes met the inclusion criteria (68KG dataset).

69 In the following, we first present results from the 29KG dataset and then evaluate the concordance of  
70 the mutational history inferred by using an expanded 68KG dataset, which establishes that the  
71 conclusions are robust to the sampling of genomes. We then applied mutational fingerprints inferred

72 using the 68KG dataset to an expanded dataset of 172,480 genomes (sampled on December 30, 2020;  
73 172KG) to track global spatiotemporal dynamics SARS-CoV-2. We have also set up a live dashboard  
74 showing regularly updated results because the processes of data analysis, manuscript preparation, and  
75 peer-review of scientific articles are much slower than the pace of expansion of SARS-CoV-2 genome  
76 collection. Also, we provide a simple “in-the-browser” tool to classify any SARS-CoV-2 genome based  
77 on key mutations derived by the MOA analysis (<http://sars2evo.datamonkey.org/>).

## 78 **Mutational history and progenitor of SARS-COV-2**

79 We used MOA to reconstruct the history of mutations that gave rise to 49 common single nucleotide  
80 variants (SNVs) in the 29KG dataset (**Fig. 1**). These variants occur with greater than 1% variant frequency  
81 ( $vf > 1\%$ ; **Fig. 2a**). For ease of reference, we used the inferred mutation history to denote key groups of  
82 mutations by assigning Greek symbols ( $\mu$ ,  $\nu$ ,  $\alpha$ ,  $\beta$ ,  $\gamma$ ,  $\delta$ , and  $\epsilon$ ) to them. Individual mutations were  
83 assigned numbers and letters based on the reconstructed order and their parent-offspring relationships  
84 (**Extended Data Table 1**). We estimated the timing of mutation for each mutation based on the  
85 timestamp of the viral samples' genome sequences in which it first appeared (**Extended Data Table 1**,  
86 see *Methods*). The inferred mutation order generally agreed with the temporal pattern of the first  
87 appearance of variants in the 29KG dataset. The sampling time of 47 out of 49 mutations was greater  
88 than or equal to the first appearance of the corresponding preceding mutation in mutational history.  
89 The exceptions were seen only for two low-frequency offshoot mutations ( $\beta_{3b}$  and  $\beta_{3c}$ ; see *Methods*).  
90 This concordance provides independent validation of the reconstructed mutation graph because  
91 neither sampling dates nor locations were used in MOA analysis.

92 We found that new variants occurred in the genomic background of the variants preceding them in the  
93 reconstructed mutation history with a very high propensity (co-occurrence index, COI > 96.7%; **Fig. 1**).  
94 This suggests a strong signal to infer a sequential mutational history. Indeed, a bootstrap analysis  
95 involving genome resampling to assess the robustness of the mutation history produced high bootstrap  
96 confidence levels (BCLs) for key groups of mutations as well as many offshoots (**Fig. 1**; BCL > 95%).  
97 However, the order of some mutations was not established with a high BCL, e.g., the relative order of  
98  $\epsilon_1$ ,  $\epsilon_2$ , and  $\epsilon_3$  mutations. This is because the three  $\epsilon$  variants almost always occur together (7,624  
99 genomes), and the intermediate combinations of  $\epsilon$  variants occurred in only 42 genomes. Similarly, the  
100 count of genomes harboring all three  $\beta$  variants (22,739 genomes) far exceeded those with two or  
101 fewer  $\beta$  variants (201 genomes). There is a strong temporal tendency of variants to be sampled together  
102 (e.g.,  $\epsilon_1 - \epsilon_3$  and  $\alpha_{1a} - \alpha_{1d}$ ), suggesting an episodic spread of variants ( $P \ll 0.01$ ; see *Methods*). This  
103 episodic spreading of variants, which do not allow for determining the precise order of mutation  
104 appearance, may be caused by founder effects, positive selection, or both (e.g., ref.<sup>28</sup>). It may  
105 sometimes be an artifact of highly uneven regional and temporal genome sequencing that will produce  
106 a biased representative sample of the actual worldwide population (**Fig. 2b**).

107 *The progenitor genome*

108 The root of the mutation tree is the most recent common ancestor (MRCA) of all the genomes analyzed,  
109 which gave rise to two early coronavirus lineages ( $\nu$  and  $\alpha$ ; **Fig. 1**). The MRCA genome was the  
110 progenitor of all SARS-CoV-2 infections globally, henceforth proCoV2, and was likely carried by the first  
111 case of human transmission in the COVID-19 pandemic (index case)<sup>20</sup>. It existed on or before December  
112 24, 2019, a date for which we have the sequence of SARS-CoV-2 infection in Wuhan, China (Wuhan-1;  
113 EPI\_ISL 402123). A comparison of proCoV2 with Wuhan-1 genomes revealed three differences in the  
114 49 positions, which was also true for other reference genomes (**Fig. 2c**). This suggests that the Wuhan-  
115 1 and the other earliest sampled genomes are derived coronavirus strains that arose from proCoV2  
116 after the divergence of  $\nu$  and  $\alpha$  lineages (**Fig. 1**). The Wuhan-1 strain evolved by three successive  $\alpha$   
117 mutations in the progenitor ( $\alpha_1$ ,  $\alpha_2$ , and  $\alpha_3$ ), a progression that is statistically supported (BCL = 100%).  
118 This high resolution is made possible by 896 intermediate genomes containing one or two  $\alpha$  variants in  
119 the 29KG dataset. Importantly, three closely-related non-human coronavirus genomes (bats and  
120 pangolin) all have the same base at these positions as does the proCoV2 genome, suggesting that the  
121 ancestral genome did not contain  $\alpha$  variants. Furthermore, genomes with  $\nu$  variants of proCoV2 do not  
122 contain the other 47 variants, all of which occurred on the genomes containing  $\alpha_1$ - $\alpha_3$  that supports the  
123 inference that coronaviruses lacking  $\alpha$  variants were the ancestors of Wuhan-1 and other genomes  
124 sampled in December 2019 in China (**Fig. 2c**). Therefore, we conclude that Wuhan-1 was not the direct  
125 ancestor of all the coronavirus infections globally.

126 Did proCoV2 propagate in the human population in 2020? A comparison of the proCoV2 genetic  
127 fingerprint (49 positions) in the 29KG collection revealed three matches in China (Fujian, Guangdong,  
128 and Hangzhou) and three in the US (Washington) in January 2020 (**Fig. 2c**). One more match was found  
129 in New York in March 2020, and the  $\nu$  mutant of proCoV2 was first sampled 59 days after the Wuhan-  
130 1 strain. This means that the progenitor coronavirus spread and mutated in the human population for  
131 weeks and months after the first reported COVID-19 cases.

132 Because proCoV2 is three bases different from the Wuhan-1 genome sampled on December 24, 2019,  
133 we estimate that the divergence of earliest variants of proCoV2 occurred 5.8 - 8.1 weeks prior based  
134 on the range of possible mutation rates of coronavirus genomes<sup>20</sup>. This timeline puts the presence of  
135 proCoV2 late-October to mid-November 2019 that is consistent with some other reports, including the  
136 report of a fragment of spike protein identical to Wuhan-1 in early December in Italy<sup>18,20,29-31</sup>. The  
137 sequenced segment of the spike protein is short (409 bases). It does not span positions in which 49  
138 major early variants were observed, which means that the Italian Spike protein fragment can only  
139 confirm the existence of proCoV2 before the first coronavirus detection in China.

140 Comparisons of the protein sequences encoded by the proCoV2 genome revealed 131 other genomic  
141 matches, which contained only synonymous differences from proCoV2. A majority (89 genomes) of  
142 these matches were from coronaviruses sampled in China and other Asian countries (**Fig. 2d**). The first  
143 sequence was sampled 12 days after the earliest sampled virus, whose genome became available on

144 December 24, 2019. Multiple matches were found in all sampled continents and detected as late as  
145 April 2020 in Europe. These spatiotemporal patterns suggest that proCoV2 already possessed the  
146 repertoire of protein sequences needed to infect, spread, and persist in the global human population  
147 (see also ref.<sup>28</sup>). Notably, none of these coronavirus genomes contained widely-studied Spike protein  
148 mutant (D614G), a  $\beta$  mutation that occurred in the genomes carrying all three  $\alpha$  variants and was first  
149 seen in late January 2020.

150 We then analyzed a later snapshot of SARS-CoV-2 genome collection, consisting of genomes obtained  
151 from GISAID, acquired three months after the 29KG dataset. This dataset expanded the collection of  
152 coronavirus genomes from viral isolates collected after July 7, 2020 (16,739 genomes) and added  
153 20,004 genome sequences from viral isolates dated before July 7, 2020. In the expanded MOA analysis,  
154 we retained 49 variants found with frequency  $> 1\%$  in the 29KG dataset and added variants found with  
155 a frequency  $> 1\%$  in the 68KG dataset (84 total variants; see **Extended Data Table 2**). MOA analysis of  
156 the 68KG dataset produced the proCoV2 genome identical to that inferred using the 29KG dataset (see  
157 *Methods*). We found one additional genome with a proCoV2 fingerprint sampled in Hubei, China, four  
158 weeks after the Wuhan-1 strain was reported.

159 The inferred mutation history from the 68KG dataset was well-supported with high COI and BCLs  
160 concordance with the mutation history produced using the 29KG dataset (**Fig. 3b**). Therefore, all the  
161 inferences reported for the 29KG dataset were robust to the expanded sampling of genomes. In the  
162 expanded mutation history, two new groups of variants were identified ( $\zeta$  and  $\eta$ ), which originated in  
163 mid-March 2020 and are found in relatively high frequency in the 68KG dataset ( $\sim 4.4\%$  and  $8.0\%$ ,  
164 respectively; **Extended Data Table 2**). Variants in  $\zeta$  and  $\eta$  groups also showed episodic accumulation of  
165 mutations, e.g., the count of genomes containing three  $\zeta$  mutations ( $\zeta_1$ - $\zeta_3$ ; 2,955 genomes) was much  
166 larger than those with a subset of these variants (148 genomes). The episodic nature of mutational  
167 spread for 84 variants in the 68KG is statistically significant ( $P < 10^{-8}$ ), i.e., clusters of mutations together  
168 have become common variants (see *Methods*).

### 169 ***Coronavirus fingerprints and spatiotemporal tracking***

170 The progression of mutations in the mutation history directly transforms into a collection of genetic  
171 fingerprints or signatures. Each fingerprint represents a genome type containing all the variants on the  
172 path from that node up to the progenitor proCoV2. These fingerprints can classify genomes and track  
173 spatiotemporal patterns of dominant lineages genomes (see *Methods*). We use a shorthand to refer to  
174 each barcode in which only the major variant type is used. For example,  $\alpha$  fingerprint refers to genomes  
175 that one or more of the  $\alpha$  variants and no other major variants, and  $\alpha\beta$  fingerprint refers to genomes  
176 that contain at least one  $\alpha$ , at least one  $\beta$  variant, and no other major variants. This nomenclature is  
177 intuitive and provides a way to glean evolutionary information from the coronavirus lineage's name. In  
178 the 68KG dataset (October 12, 2020 GISAID snapshot), global frequencies of major proCoV2 fingerprints  
179 were  $\alpha\beta\epsilon$  (32.1%),  $\alpha\beta\gamma\delta$  (17.7%),  $\alpha\beta$  (16.7%),  $\alpha\beta\eta$  (9.9%),  $\alpha\beta$  (9.8%),  $\alpha\beta\gamma$  (6.8%),  $\alpha\beta\zeta$  (4.5%), and  $\nu$   
180 (2.4%).

181 **Figure 4** shows the evolving spatiotemporal of all major fingerprints in Asia, Europe, and North America  
182 inferred for an expanded dataset of 172,480 genomes (December 30, 2020 snapshot). Spatiotemporal  
183 patterns in cities, countries, and other regions are available online at <http://sars2evo.datamonkey.org/>.  
184 We observe the spread and replacement of prevailing strains in Europe ( $\alpha\beta\epsilon$  with  $\alpha\beta\zeta$ ) and Asia ( $\alpha$  with  
185  $\alpha\beta\epsilon$ ), the preponderance of the same strain for most of the pandemic in North America ( $\alpha\beta\gamma\delta$ ), and the  
186 continued presence of multiple high-frequency strains in Asia and North America. Spatiotemporal  
187 patterns of strain spread converged for Europe and Asia by July-August 2020 to  $\alpha\beta\epsilon$  genetic fingerprints.  
188 These patterns diverged from North America, where  $\alpha\beta$  along with its mutant ( $\alpha\beta\gamma\delta$ ) were common.  
189 After that, Europe saw  $\zeta$  variants of  $\alpha\beta$  grow ( $\alpha\beta\zeta$ ), replacing  $\alpha\beta\epsilon$  genomes and its new  $\eta$  offshoot ( $\alpha\beta\epsilon\eta$ )  
190 (e.g., ref.<sup>32</sup>). The  $\zeta$  mutations were first detected three weeks after the sampling of the first  $\epsilon$  variants.  
191 Remarkably,  $\alpha\beta\gamma\delta$  has remained the dominant lineage in North America since April 2020, in contrast to  
192 the turn-over seen in Europe and Asia. More recently, novel fast-spreading variants have been reported  
193 (e.g., ref<sup>33</sup>). In particular, an S protein variant (N501Y) from South Africa and London has rapidly  
194 increased<sup>33</sup>. Coronaviruses with N501Y variant in South Africa carry the  $\alpha\beta\gamma\delta$  genetic fingerprint,  
195 whereas those in London carry the  $\alpha\beta\epsilon$  genetic fingerprint. This means that the N501Y mutation arose  
196 independently in two coronavirus lineages that show convergent patterns of increased spread. At  
197 present,  $\alpha\beta\zeta$  dominates the UK, and the number of genomes publicly available from South Africa is  
198 relatively small to make reliable inferences at present (see <http://sars2evo.datamonkey.org> for future  
199 updates). Overall, our mutational fingerprinting and nomenclature provides a simple way to glean the  
200 ancestry of new variants in contrast to phylogenetic designations (e.g., B.1.350 and B.1.1.7<sup>33</sup>).

## 201 **Conclusions**

202 Through innovative analyses of two large collections of SARS-CoV-2 genomes, we have consistently  
203 reconstructed the same progenitor coronavirus genome and identified its presence worldwide for  
204 many months after the pandemic began. The progenitor genome is a better reference for rooting  
205 phylogenies, orienting mutations, and estimating sequence divergences. The reconstructed mutational  
206 history of SARS-CoV-2 revealed major mutational fingerprints to identify and track the novel  
207 coronavirus's spatiotemporal evolution, revealing convergences and divergences of dominant strains  
208 among geographical regions from an analysis of more than 174 thousand genomes.

209 Furthermore, the approach taken here to reconstruct the progenitor genome and discover key  
210 mutational events will generally be applicable for analyzing pathogens during the early stages of  
211 outbreaks. The approach is scalable for even bigger datasets because it does not require more  
212 phylogenetically informative variants with an increasing number of samples. In fact, it benefits from  
213 bigger datasets as they afford more accurate estimates of individual and co-occurrence frequencies of  
214 variants and enable more reliable detection of lower frequency variants. Its continued application to  
215 SARS-CoV-2 genomes and other pathogen outbreaks will produce their ancestral genomes and their

216 spatiotemporal dynamics, improving our understanding of the past, current, and future evolution of  
217 pathogens and associated diseases.

## 218 **Methods**

### 219 Genome data acquisition and processing

220 We first downloaded 60,332 SARS-CoV-2 genomes from the GISAID<sup>27</sup> database, along with information  
221 on sample collection dates and locations (until July 7, 2020). Of all the genomes downloaded, we only  
222 retained those with greater than 28,000 bases and were marked as originating from human hosts and  
223 passing controls detailed below. Similarly, the second dataset, the 68KG dataset, was assembled from  
224 133,741 genomes and downloaded on October 12, 2020. Again, we retained only those with greater  
225 than 28,000 bases and marked as originating from human hosts.

226 Each genome was subjected to codon-aware alignment with the NCBI reference genome (accession  
227 number NC\_045512) and then subdivided into ten regions based on CDS features: ORF1a (including  
228 nsp10), ORF1b (starting with nsp12), S, ORF3a, E, M, ORF6, ORF7a, ORF8, N, and ORF10. Gene ORF7b  
229 was removed because it was too short for alignment and comparisons. For each region, we scanned  
230 and discarded sequences containing too many ambiguous nucleotides to remove data with too many  
231 sequencing errors. Thresholds were 0.5% for the S gene, 0.1% for ORF1a and ORF1b genes, and 1% for  
232 all other genes. We mapped individual sequences to the NCBI reference genome (NC\_045512) using a  
233 codon-aware extension to the Smith-Waterman algorithm implemented in HyPhy<sup>34</sup>  
234 (<https://github.com/veg/hyphy-analyses/tree/master/codon-msa>), translated mapped sequence to  
235 amino-acids, and performed multiple protein sequence alignment with the auto settings function of  
236 MAFFT (version 7.453)<sup>35</sup>. Codon sequences were next mapped onto the amino-acid alignment. The  
237 multiple sequence alignment of SARS-CoV-2 genomes was aligned with the sequence of three closest  
238 outgroups, including the coronavirus genomes of the *Rhinolophus affinis* bat (RaTG13), *R. malayanus*  
239 bat (RmYN02), and *Manis javanica* pangolin (MT121216.1)<sup>36,37</sup>. The alignment was visually inspected  
240 and adjusted in Geneious Prime 2020.2.2 (<https://www.geneious.com>). The final alignment contained  
241 all genomic regions except ORF7b and non-coding regions (5' and 3' UTRs, and intergenic spacers). After  
242 these filtering and alignment steps, the multiple sequence alignment contained 29,115 sites and 29,681  
243 SARS-CoV-2 genomes for the July 7, 2020 snapshot, which we refer to as the 29KG dataset. For the  
244 October 12 snapshot, there were 68,057 sequences, which we refer to as the 68KG dataset. We also  
245 conducted a spatiotemporal analysis on an expanded dataset containing 172,480 genomes (172KG)  
246 acquired on December 30, 2020.

### 247 Reference genomes and collection dates

248 We used the dates of viral collections provided by the GISAID database<sup>27</sup> in all our analyses if they were  
249 resolved to the day (i.e., we discarded data that only contained partial dates, e.g., April 2020). All  
250 genomes were used in the mutation ordering analyses, but genomes with incomplete sampling dates



251 were excluded from the spatiotemporal analyses and derived interpretations. We noted that the  
252 earliest sample included in GISAID (ID: EPI\_ISL\_402123) was collected on December 24, 2019, although  
253 the NCBI website lists its collection date as December 23, 2019 (GenBank ID: MT019529). Therefore,  
254 we used the GISAID collection date for the sake of consistency. Regarding the NCBI reference genome  
255 (GenBank ID: NC\_045512; GISAID ID: EPI\_ISL\_402125)<sup>38</sup>, this sample was collected on December 26,  
256 2019<sup>39</sup>. We also used the GIS reference genome in our analysis (ID: EPI\_ISL\_402124), collected on  
257 December 30, 2019<sup>40</sup>.

#### 258 Mutation order analyses (MOA)

259 First, we analyzed the 29KG dataset. We used a maximum likelihood method, SCITE<sup>23</sup>, and variant co-  
260 occurrence analyses for reconstructing the order of mutations corresponding to 49 common variants  
261 (frequency > 1%) observed in this dataset. MOA has demonstrated high accuracy for analyzing tumor  
262 cell genomes that reproduce clonally, have frequent sequencing errors, and exhibit limited sequence  
263 divergence<sup>23,24</sup>. In MOA, higher frequency variants are expected to have arisen earlier than low-  
264 frequency variants in clonally reproducing populations<sup>23,26</sup>. We used the highest frequency variants to  
265 anchor the analysis and the shared co-occurrence of variants among genomes to order mutations while  
266 allowing probabilistically for sequencing errors and pooled sequencing of genomes<sup>23</sup>. MOA is different  
267 from traditional phylogenetic approaches where positions are treated independently, i.e., the shared  
268 co-occurrence of variants is not directly utilized in the inference procedure. Notably, both traditional  
269 phylogenetic and mutation order analyses are expected to produce concordant patterns when  
270 sequencing errors and other artifacts are minimized. However, sequencing errors and limited  
271 mutational input during the coronavirus history adversely impact traditional methods, as does the fact  
272 that the closest coronaviruses useable as outgroups have more than a thousand base differences from  
273 SARS-CoV-2 genomes that only differ in a handful of bases from each other<sup>7,9,10</sup>.

274 MOA requires a binary matrix of presence/absence (1/0) of mutants, which is straightforward in  
275 analyzing cell sequences from tumors because they arise from normal cells that supply the definitive  
276 ancestral state. To designate mutation orientations for applying MOA in SARS-CoV-2 analysis, we  
277 devised a simple approach in which we began by comparing nucleotides at the 49 genomic positions  
278 among three closely-related genomes (bat RaTG13, bat RmYN02, and pangolin MT121216.1)<sup>41</sup>. We  
279 chose the consensus base to be the initial reference base, such that SARS-CoV-2 genome bases were  
280 coded to be "0" whenever they were the same as the consensus base at their respective positions. All  
281 other bases were assigned a "1." There were 39 positions in which all three outgroup genomes were  
282 identical to each other and 9 in which two of the outgroups showed the same base. In the remaining  
283 position (28657), all three outgroups differed, so we selected the base found in the gene with the  
284 highest sequence similarity to the human SARS-CoV-2 NCBI reference genome (NC\_045512) because  
285 SARS-CoV-2's ancestor likely experienced genomic recombination before its zoonotic transfer into  
286 humans<sup>28,42,43</sup>. At one position, both major and minor bases in humans were different from the

287 consensus base in the outgroups, so we assigned the mutant status to the minority base ( $U$ ;  $vf = 29.8\%$ ).  
288 All missing and ambiguous bases were coded to be ignored (missing data) in all the analyses.

289 These initially assigned mutation orientations were tested in a subsequent investigation of variants' co-  
290 occurrence index (COI). COI for a given variant ( $y$ ) is the number of genomes that contain  $y$  and its  
291 directly preceding mutation ( $x$ ) in the mutation history, divided by the number of genomes that contain  
292  $y$ . When COI was lower than 70%, we reversed each position's mutation orientation individually and  
293 selected the mutation orientation that produced mutation histories with the highest COI.

294 In the SCITE analysis of 49 variants and 29,861 genomes, we started with default parameter settings of  
295 false-negative rate (FNR = 0.21545) and false-positive rate (FPR = 0.0000604) of mutation detection.  
296 We carried out five independent runs to ensure stability and convergence to obtain 29KG collection-  
297 specific estimates of FNR and FPR by comparing the observed and predicted sequences based on this  
298 mutation graph. The estimated FNR (0.00488) and FPR (0.00800) were very different from the SCITE  
299 default parameters, where the estimated FNR was much lower. This difference in error rates is  
300 unsurprising because we used only common variants ( $vf > 1\%$ ), and the 29KG dataset was not obtained  
301 from single-cell sequencing in which dropout during single-cell tumor sequencing elevates FNR, i.e.,  
302 mutant alleles are not sequenced.

303 As noted above, the initial mutation orientations were simply the starting designations for our analysis,  
304 which are subsequently investigated by evaluating the COI of each variant in the reconstructed  
305 mutation history. In this process, we reverse ancestor/mutant coding for variants that received low COI  
306 to examine if a mutation history with higher COI can be generated. Two positions (3037 and 28854)  
307 received low COI ( $< 70\%$ ). At position 3037, the reversed encoding ( $C \rightarrow U$ ) received significantly higher  
308 COI (100%) than the starting encoding ( $U \rightarrow C$ ; 60%), so the position was recoded. At position 28854,  
309 the ordering and direction of mutation remained ambiguous despite extensive analyses, but it did not  
310 impact the predicted MRCA sequence. Therefore, we only recoded the column for position 3037 and  
311 generated a new  $49 \times 29861$  (SNVs  $\times$  genomes) matrix to conduct a SCITE analysis.

312 At one position (28657), all three outgroup sequences had different bases, so we initially selected the  
313 base found in the gene with the highest sequence similarity to the human SARS-CoV-2 NCBI reference  
314 genome. We next tested if reversed encoding produced a better mutation graph. The reversed  
315 encoding produced a mutation graph with a much higher log-likelihood ( $-32355.58$  and  $-30289.92$ , for  
316 the initial and reversed encoding, respectively;  $P \ll 0.01$  using the AIC protocol in ref.<sup>44</sup>). Therefore,  
317 we recoded position 28657 and generated a new  $49 \times 29861$  (SNVs  $\times$  genomes) matrix.

318 It was subjected to SCITE analysis and produced a mutation graph for 49 variants in the 29KG dataset.  
319 This graph predicts an FNR of 0.00418 and FPR of 0.00295 per base. Using these new FNR and FPR, we  
320 again performed SCITE analysis and produced the final mutation history graph. Starting from the top of  
321 a mutation graph, a distinct Greek symbol was assigned to a group of mutations that were occurred  
322 sequentially, and variants with similar frequency were assigned the same Greek symbol ( $\mu$ ,  $\nu$ ,  $\alpha$ ,  $\beta$ ,  $\gamma$ ,  $\delta$ ,  
323 and  $\epsilon$ ). The high-frequency variants with the same Greek symbol were distinguished by numbers to

324 represent the sequential relationship, e.g.,  $\alpha_1$  and  $\alpha_2$ . When an offshoot of a high-frequency mutation  
325 had low variant frequency, we assigned it the same Greek symbol and number to represent the parent-  
326 offspring relationship and further distinguished descendants by adding a small letter, e.g.,  $\alpha_{1a}$  and  $\alpha_{1b}$ .

327 In this mutation graph, the most recent common ancestor (MRCA) corresponds to the progenitor that  
328 gave rise to  $v$  and  $\alpha$  lineages. MRCA is the progenitor of all human SARS-CoV-2 infections (proCoV2),  
329 which descended from the parental lineage that divergence form and its closest relatives, including bats  
330 and pangolins. We estimate that proCoV2 existed 5.8 to 8.1 weeks before December 24, 2019, on which  
331 the Wuhan-1 was sampled, by using SARS-CoV-2 HPD mutation rate range of  $6.64 \times 10^{-4} - 9.27 \times 10^{-4}$   
332 substitutions per site per year<sup>20</sup>. We have made available the proCoV2 genome sequence in FastA  
333 format at <http://igem.temple.edu/COVID-19>, which is the same as the NCBI reference genome with  
334 base differences corresponding to  $\alpha_1 - \alpha_3$  mutations at positions 18060, 8782, and 28144, as discussed  
335 in the main text. In this mutation graph, COI for each variant is shown next to the arrow.

### 336 Bootstrap analysis

337 We assessed the robustness of the mutation history inference to genome sampling by bootstrap  
338 analysis. We generated 100 bootstrap replicate datasets, each built by randomly selecting 29,861  
339 genomes with replacement. Then, SCITE was used to infer the mutation graph for each replicate  
340 dataset. Bootstrap confidence level, scored for each variant pair, was the number of replicates in which  
341 the given pair of variants were directly connected in the mutation history in the same way as shown in  
342 **figure 1**. BCLs were often lower for major variants within groups (e.g.,  $\epsilon_1 - \epsilon_3$ ) because they occur with  
343 very similar frequencies. This feature adversely affected the BCL values of mutation orders between  
344 groups, e.g.,  $\beta$  and  $\epsilon$ . In this case, we considered each group as a single entity. We computed BCL to be  
345 the proportion of replicates in which pairs of groups were directly connected in the mutation history in  
346 the same way as shown in **figure 1**. Groups used were  $\beta_1 - \beta_3$ ,  $\epsilon_1 - \epsilon_3$ , and  $\alpha_{1a} - \alpha_{1d}$ . All of these BCL values  
347 are shown with an underline.

### 348 Temporal concordance

349 Because mutation ordering analysis analyses did not use spatial or temporal information for genomes  
350 or mutations, it can be validated by evaluating the concordance of the inferred order of mutations with  
351 the timing of their first appearance ( $tf$ ). Using the genomes for which virus sampling day, month, and  
352 year were available, we determined  $tf$  for every variant in the 29KG dataset. For a mutation  $i$ , we  
353 compared its  $tf(i)$  with  $tf(j)$  such that  $j$  is the nearest preceding mutation in the mutation graph. We  
354 found that  $tf(j) \geq tf(i)$  for 47 of 49 mutations, except for  $\beta_{3b}$  and  $\beta_{3c}$  pairs. These two offshoot mutants  
355 of  $\beta_3$  were sampled 35 days ( $\beta_{3b}$ ) and 12 days ( $\beta_{3c}$ ) earlier than their predecessors, which could be due  
356 to their low frequency or sequencing error. COI of one variant ( $\beta_{3b}$ ) was low (54%), but the other variant  
357 ( $\beta_{3c}$ ) had a very high COI (97%).

### 358 Mutational fingerprints

359 Each node in the mutational history graph predicts an intermediate (ancestral) or a tip sequence,  
360 containing all the mutations from that node to the mutation graph's root. The mutational fingerprint is  
361 then produced directly from the mutation history graph drawn as a directional graph anchored on the  
362 root node. We compared our mutational fingerprints of the genomes in the 29KG dataset with a  
363 phylogeny-based classification<sup>1</sup> obtained using the Pangolin service (v2.0.3; [https://pangolin.cog-  
364 uk.io/](https://pangolin.cog-uk.io/)). We assigned each of the 29K genomes to a fingerprint based on the highest sequence similarity  
365 at positions containing 49 common variants. Mismatches were allowed, as sequencing errors could  
366 create them. A small fraction of genomes (1.8%) could not be assigned unambiguously to one  
367 fingerprint, so they were excluded and investigated in the future. The number of genomes assigned to  
368 each fingerprint is shown in **Extended Data Table 1**. We submitted genome sequences to the Pangolin  
369 website for classification one-by-one, and a clade designation was received. The results are summarized  
370 in **Extended Data Figure 1**. In this table, all phylogenetic-groups with fewer than 20 genomes were  
371 excluded.

372 Of the 80 phylogenetic groups shown, 74 are defined primarily by a single mutation-based fingerprint,  
373 as more than 90% of the genomes in those phylogenetic groups share the same fingerprint. This  
374 includes all small and medium-sized phylogenetic groups (up to 488 genomes) and two large groups  
375 (A.1 with 1,377 genomes and B.1.2 with 749 genomes). One large group, B.1.1, predominately connects  
376 with  $\epsilon_3$  node (79%, 4,832 genomes), but some of its members belong to  $\epsilon_3$  offshoots because they  
377 contain respective diagnostic mutations. For group B.1.1.1, two other  $\epsilon_3$  offshoots are mixed up almost  
378 equally. Three other large differences between mutational fingerprint-based classification and  
379 phylogeny-based grouping are seen for A, B, B1.1, and B.2 groups. These differences are likely because  
380 the location of the root and the earliest branching order of coronavirus lineages are problematic in  
381 phylogeny-based classifications<sup>7,9,10,14</sup>. Overall, our mutational fingerprints are immediately informative  
382 about the mutational ancestry of genomes.

### 383 Analysis of 68KG dataset

384 We repeated the above MOA procedure on the 68KG dataset (68,057 genomes). This 68KG data  
385 contained 72 common variants (>1% frequency). For direct comparison purposes, we added 12 variants  
386 that were common variants on 29KG data, but their frequency had become less than 1% in the 68KG  
387 data. Therefore, we used 84 variants in total and constructed a matrix of 84 × 68,057 (SNVs × genomes)  
388 for the SCITE analysis to determine the mutational order. We also conducted the bootstrap analysis  
389 and assigned mutational fingerprints using the procedure mentioned above. The number of genomes  
390 mapped to each fingerprint is listed in **Extended Data Table 2**.

### 391 Spatiotemporal analysis of 172KG dataset

392 We developed a sequence classification protocol that first calls variants in a viral genome using proCoV2  
393 as the reference sequence using minimap2<sup>45</sup>. Then, it assigns the sequence to a path in the mutation  
394 graph with the highest concordance (Jaccard index). It is implemented in a simple browser-based tool,  
395 which shows the example output for ENA accession number MT675945 (**Extended Data Figure 2**;

396 <http://sars2evo.datamonkey.org>). The classification is conducted on the client-side such that the  
397 researcher's data never leaves their personal computer.

#### 398 Testing episodic spread of variants

399 We performed non-parametric Wald–Wolfowitz runs-tests<sup>46,47</sup> of the null hypothesis that the first  
400 sampling of variants is randomly distributed over time (i.e., evenly spaced). The null hypothesis was  
401 rejected for both 29KG and 64KG analysis at  $P \ll 0.01$ , suggesting significant temporally clustering in  
402 both 29KG dataset and 64KG datasets. Because many mutations were first sampled on December 24,  
403 2019, we only included one mutation for that day to avoid biasing the test.

404 **Data Availability and Code Availability:** Live evolutionary history and spatiotemporal distributions of  
405 common variants can be accessed via <http://igem.temple.edu/COVID-19> (beta version). All genome  
406 sequences and metadata are available publicly at GISAID (<https://www.gisaid.org/>), and the predicted  
407 proCoV2 sequence is available at <http://igem.temple.edu/COVID-19>. The other relevant information is  
408 provided in the supplementary materials.

409

#### 410 **References**

- 411 1. Rambaut, A. *et al.* A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic  
412 epidemiology. *Nat. Microbiol.* (2020) doi:10.1038/s41564-020-0770-5.
- 413 2. Worobey, M. *et al.* The emergence of SARS-CoV-2 in Europe and the US. *bioRxiv* (2020)  
414 doi:10.1101/2020.05.21.109322.
- 415 3. Tang, X. *et al.* On the origin and continuing evolution of SARS-CoV-2. *Natl. Sci. Rev.* **7**, 1012–1023  
416 (2020).
- 417 4. Lemey, P. *et al.* Accommodating individual travel history and unsampled diversity in Bayesian  
418 phylogeographic inference of SARS-CoV-2. *Nat. Commun.* **11**, 1–14 (2020).
- 419 5. Forster, P., Forster, L., Renfrew, C. & Forster, M. Phylogenetic network analysis of SARS-CoV-2  
420 genomes. *Proc. Natl. Acad. Sci. U. S. A.* **117**, 9241–9243 (2020).
- 421 6. Turakhia, Y. *et al.* Stability of SARS-CoV-2 phylogenies. *PLOS Genetics* vol. 16 (2020).
- 422 7. Mavian, C. *et al.* Sampling bias and incorrect rooting make phylogenetic network tracing of SARS-  
423 COV-2 infections unreliable. *Proc. Natl. Acad. Sci. U. S. A.* **117**, 12522–12523 (2020).
- 424 8. De Maio, N. *et al.* Issues with SARS-CoV-2 sequencing data. [https://virological.org/t/issues-with-](https://virological.org/t/issues-with-sars-cov-2-sequencing-data/473)  
425 [sars-cov-2-sequencing-data/473](https://virological.org/t/issues-with-sars-cov-2-sequencing-data/473).
- 426 9. Morel, B. *et al.* Phylogenetic analysis of SARS-CoV-2 data is difficult. *bioRxiv* (2020)  
427 doi:10.1101/2020.08.05.239046.
- 428 10. Pipes, L., Wang, H., Huelsenbeck, J. & Nielsen, R. Assessing uncertainty in the rooting of the  
429 SARS-CoV-2 phylogeny. *Mol. Biol. Evol.* (2020) doi:10.1101/2020.06.19.160630.
- 430 11. Lai, A., Bergna, A., Acciarri, C., Galli, M. & Zehender, G. Early phylogenetic estimate of the  
431 effective reproduction number of SARS-CoV-2. *J. Med. Virol.* **92**, 675–679 (2020).
- 432 12. Castells, M., Lopez-Tort, F., Colina, R. & Cristina, J. Evidence of Increasing Diversification of  
433 Emerging SARS-CoV-2 Strains. *J. Med. Virol.* 1–8 (2020).
- 434 13. Andersen, K. G., Rambaut, A., Lipkin, W. I., Holmes, E. C. & Garry, R. F. The proximal origin of  
435 SARS-CoV-2. *Nat. Med.* **26**, 450–452 (2020).
- 436 14. Wenzel, J. Origins of SARS-CoV-1 and SARS-CoV-2 are often poorly explored in leading  
437 publications. *Cladistics* **36**, 374–379 (2020).
- 438 15. Gómez-carballa, A., Bello, X., Pardo-seco, J., Martinon-Torres, F. & Salas, A. genome variation  
439 worldwide highlights the impact of COVID-19 super-spreaders. *Genome Res.* (2020)  
440 doi:10.1101/gr.266221.120.
- 441 16. Fauver, J. R. *et al.* Coast-to-Coast Spread of SARS-CoV-2 during the Early Epidemic in the United  
442 States. *Cell* **181**, 990–996.e5 (2020).

- 443 17. Dearlove, B. L. *et al.* A SARS-CoV-2 vaccine candidate would likely match all currently circulating  
444 strains. *bioRxiv* (2020) doi:10.1101/2020.04.27.064774.
- 445 18. Stefanelli, P. *et al.* Whole genome and phylogenetic analysis of two SARS-CoV-2 strains isolated  
446 in Italy in January and February 2020: Additional clues on multiple introductions and further  
447 circulation in Europe. *Eurosurveillance* **25**, 1–5 (2020).
- 448 19. Lu, R. *et al.* Genomic characterisation and epidemiology of 2019 novel coronavirus: implications  
449 for virus origins and receptor binding. *Lancet* **395**, 565–574 (2020).
- 450 20. Pekar, J., Worobey, M., Moshiri, N., Scheffler, K. & Wertheim, J. O. Timing the SARS-CoV-2 Index  
451 Case in Hubei Province. *bioRxiv* 2020.11.20.392126 (2020).
- 452 21. Nei, M. & Kumar, S. *Molecular Evolution and Phylogenetics*. (Oxford University Press, 2002).
- 453 22. Yang, Z., Kumar, S. & Nei, M. A new method of inference of ancestral nucleotide and amino acid  
454 sequences. *Genetics* **141**, 1641–1650 (1995).
- 455 23. Jahn, K., Kuipers, J. & Beerenwinkel, N. Tree inference for single-cell data. *Genome Biol.* **17**, 1–  
456 17 (2016).
- 457 24. Miura, S. *et al.* Computational enhancement of single-cell sequences for inferring tumor  
458 evolution. *Bioinformatics* **34**, i917–i926 (2018).
- 459 25. Ross, E. M. & Markowitz, F. OncoNEM: Inferring tumor evolution from single-cell sequencing  
460 data. *Genome Biol.* **17**, 1–14 (2016).
- 461 26. Kim, K. I. & Simon, R. Using single cell sequencing data to model the evolutionary history of a  
462 tumor. *BMC Bioinformatics* **15**, (2014).
- 463 27. Hadfield, J. *et al.* NextStrain: Real-time tracking of pathogen evolution. *Bioinformatics* **34**, 4121–  
464 4123 (2018).
- 465 28. MacLean, O. A. *et al.* Natural selection in the evolution of SARS-CoV-2 in bats, not humans,  
466 created a highly capable human pathogen. *bioRxiv* (2020) doi:10.1101/2020.05.28.122366.
- 467 29. van Dorp, L. *et al.* Emergence of genomic diversity and recurrent mutations in SARS-CoV-2.  
468 *Infect. Genet. Evol.* **83**, 104351 (2020).
- 469 30. Li, X. *et al.* Transmission dynamics and evolutionary history of 2019-nCoV. *J. Med. Virol.* **92**, 501–  
470 511 (2020).
- 471 31. Giovanetti, M., Benvenuto, D., Angeletti, S. & Ciccozzi, M. The first two cases of 2019-nCoV in  
472 Italy: Where they come from? *J. Med. Virol.* **92**, 518–521 (2020).
- 473 32. Hodcroft, E. B. *et al.* Emergence and spread of a SARS-CoV-2 variant through Europe in the  
474 summer of 2020. *medRxiv* 2020.10.25.20219063 (2020).
- 475 33. Rambaut, A. *et al.* Preliminary genomic characterisation of an emergent SARS-CoV-2 lineage in  
476 the UK defined by a novel set of spike mutations - SARS-CoV-2 coronavirus / nCoV-2019 Genomic  
477 Epidemiology - Virological. 2020 [https://virological.org/t/preliminary-genomic-characterisation-](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)  
478 [of-an-emergent-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).
- 479 34. Gianella, S. *et al.* Detection of Minority Resistance during Early HIV-1 Infection: Natural Variation  
480 and Spurious Detection rather than Transmission and Evolution of Multiple Viral Variants. *J.*  
481 *Virology* **85**, 8359–8367 (2011).
- 482 35. Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7:  
483 Improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772–780 (2013).
- 484 36. Zhou, P. *et al.* A pneumonia outbreak associated with a new coronavirus of probable bat origin.  
485 *Nature* **579**, 270–273 (2020).
- 486 37. Liu, P. *et al.* Are pangolins the intermediate host of the 2019 novel coronavirus (SARS-CoV-2)?  
487 *PLoS Pathog.* **16**, 1–13 (2020).
- 488 38. Wu, F. *et al.* A new coronavirus associated with human respiratory disease in China. *Nature* **579**,  
489 265–269 (2020).
- 490 39. Chiara, M., Horner, D. S., Gissi, C. & Pesole, G. Comparative genomics provides an operational  
491 classification system and reveals early emergence and biased spatio-temporal distribution of  
492 SARS-CoV-2. *bioRxiv* (2020) doi:10.1101/2020.06.26.172924.
- 493 40. Okada, P. *et al.* Early transmission patterns of coronavirus disease 2019 (COVID-19) in travellers  
494 from Wuhan to Thailand, January 2020. *Eurosurveillance* **25**, 2000097 (2020).
- 495 41. Boni, M. F. *et al.* Evolutionary origins of the SARS-CoV-2 sarbecovirus lineage responsible for the  
496 COVID-19 pandemic. *Nat. Microbiol.* (2020) doi:10.1038/s41564-020-0771-4.
- 497 42. Li, X. *et al.* Emergence of SARS-CoV-2 through recombination and strong purifying selection. *Sci.*  
498 *Adv.* **6**, 1–12 (2020).
- 499 43. Huang, J.-M., Jan, S. S., Wei, X., Wan, Y. & Ouyang, S. Evidence of the Recombinant Origin and

- 500 Ongoing Mutations in Severe Acute Respiratory Syndrome Coronavirus 2 (SARS-CoV-2). *bioRxiv*  
501 (2020) doi:10.1101/2020.03.16.993816.  
502 44. Pupko, T., Huchon, D., Cao, Y., Okada, N. & Hasegawa, M. Combining multiple data sets in a  
503 likelihood analysis: Which models are the best? *Mol. Biol. Evol.* **19**, 2294–2307 (2002).  
504 45. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100  
505 (2018).  
506 46. Wald, A. & Wolfowitz, J. On a Test Whether Two Samples are from the Same Population. *Ann.*  
507 *Math. Stat.* **11**, 147–162 (1940).  
508 47. Mateus, A. & Caeiro, F. An R implementation of several randomness tests. *AIP Conf. Proc.* **1618**,  
509 531–534 (2015).  
510

## 511 **Acknowledgments**

512 We thank all the authors and organizations who have kindly deposited and shared genome data on  
513 GISAID (see <http://igem.temple.edu/COVID-19> for a list of all the authors). We thank Ananias Escalante,  
514 Rob Kulathinal, Li Liu, Jose Barba-Montoya, Antonia Chroni, Ravi Patel, and Caryn Babaian for their  
515 critical comments. We appreciate the technical support provided by Jared Knoblauch and Glen Stecher.  
516 This research was supported by grants from the U.S. National Science Foundation to S.K. (GCR-1934848,  
517 DEB-2034228) and S.P. (DBI-2027196) and from the U.S. National Institutes of Health to S.K. (GM-  
518 0126567-03 and 139504-01) and S.P. (AI-134384).

## 519 **Author Contributions**

520 S.K. and S.M. conceived the project, designed analyses and visualizations, conducted initial analyses,  
521 and wrote the manuscript. S.P., S.W., and S.K. designed and developed the browser resource and tools.  
522 S.P., S.W., and M.A.C.O. assembled sequence alignments. M.A.C.O., S.M., S.S., and Q.T. conducted  
523 analyses and rendered visualizations. All authors intellectually contributed by discussing results and  
524 patterns, and everyone contributed to writing the manuscript.

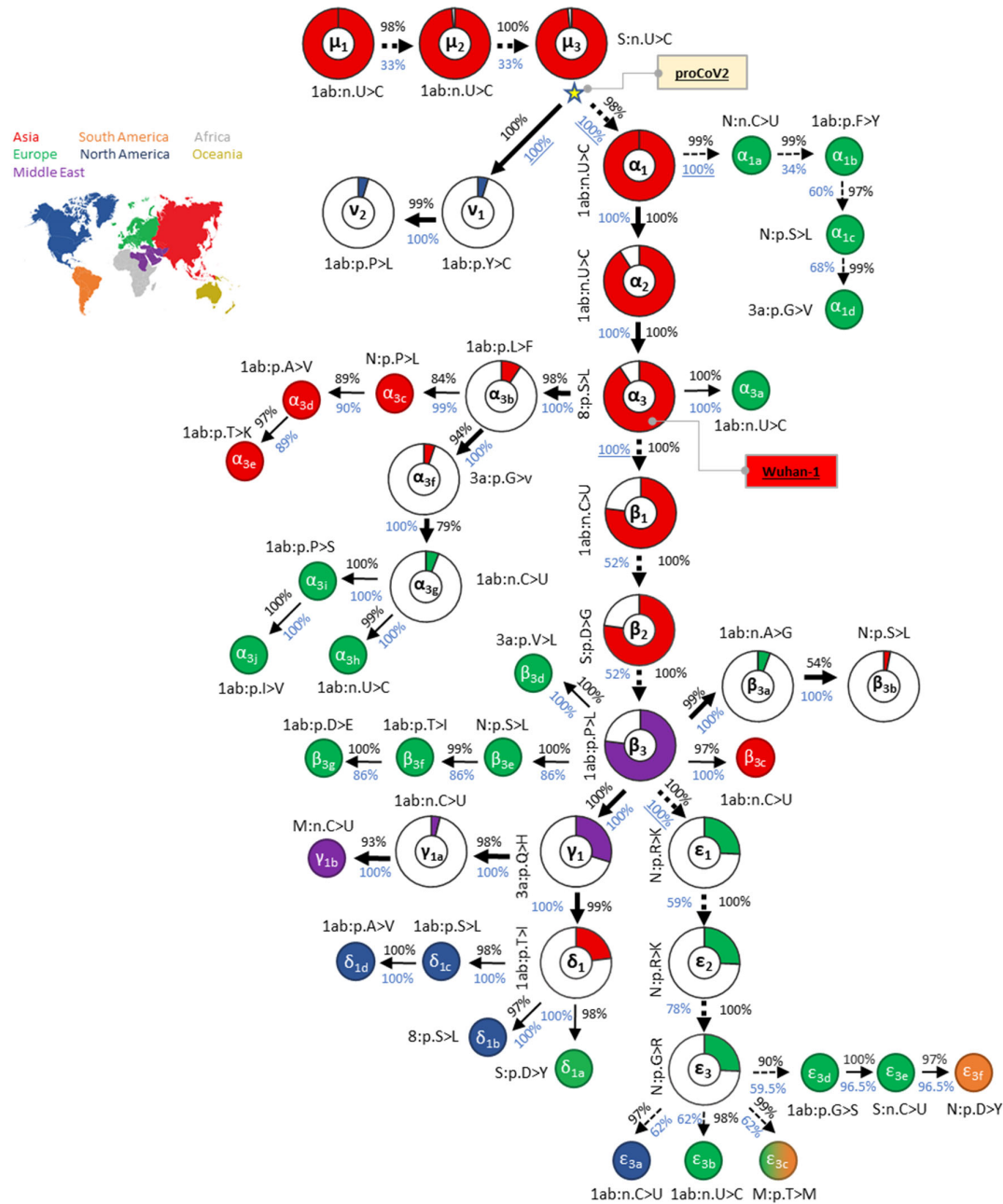
## 525 **Competing Interests**

526 The authors declare that they have no competing interests.

## 527 **Additional Information**

528 **Supplementary Information** is available for this paper. Correspondence and requests for materials  
529 should be addressed to [s.kumar@temple.edu](mailto:s.kumar@temple.edu).

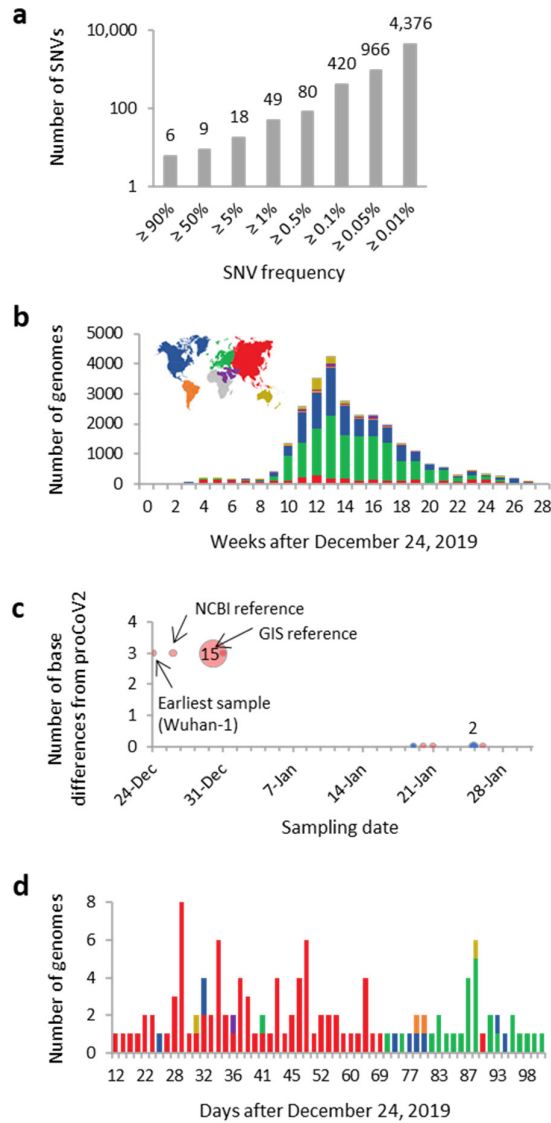
530



531

532 **Fig. 1.** Mutational history graph of SARS-CoV-2 from the 29KG dataset. Thick arrows mark the pathway of  
533 widespread variants (frequency,  $vf \geq 3\%$ ), and thin arrows show paths leading to other common mutations ( $3\% >$   
534  $vf > 1\%$ ). The pie-charts' size is proportional to variant frequency in the 29KG dataset, with pie-charts shown for  
535 variants with  $vf > 3\%$  and pie color based on the world's region where that mutation was first observed. A circle is  
536 used for all other variants, with the filled color corresponding to the earliest sampling region. The co-occurrence  
537 index (COI, black font) and the bootstrap confidence level (BCL, blue font) of each mutation and its predecessor  
538 mutation are shown next to the arrow connecting them. Underlined BCL values mark variant pairs for which BCLs  
539 were estimated for groups of variants (<80%; dashed arrows). Base changes (n.) are shown for synonymous mutations,  
540 and amino acid changes (p.) are shown for nonsynonymous mutations along with the gene/protein names ("ORF"  
541 is omitted from gene name abbreviations given in **Extended Data Table 1**). More details on each mutation are  
542 presented in **Extended Data Table 1**.  
543

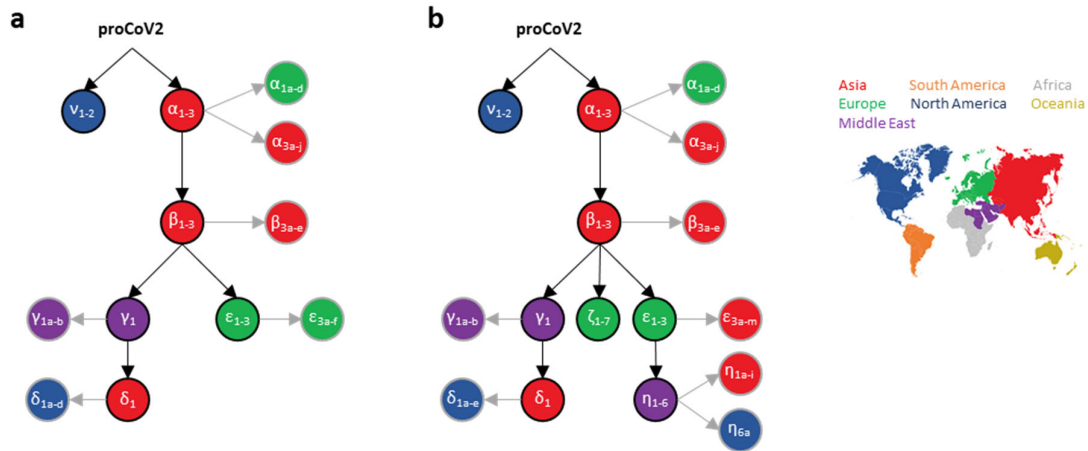




544

545 **Fig. 2. Counts of single nucleotide variants (SNVs) and genomes in the 29KG dataset.** (a) Cumulative count  
 546 of SNVs presented in the 29KG genome dataset at different frequencies. (b) The number of genomes in the 29KG  
 547 collection that were isolated weekly during the pandemic. (c) The number of base differences from proCoV2 for  
 548 genomes that were sampled in December 2019 and January 2020. The 18 genomes sampled in December 2019  
 549 in China (red) have three common SNVs different from proCoV2. In contrast, six genomes sampled in January  
 550 2020 in China (Asia, red) and the US (North America, blue) show no base differences. Multiple genomes (2 and  
 551 15) were sampled on two different days. (d) Temporal and spatial distribution of strains identical to proCoV2 at the  
 552 protein sequence level, i.e., they have only  $\mu$  mutations. The color scheme used to mark sampling locations is  
 553 shown in panel b.

554

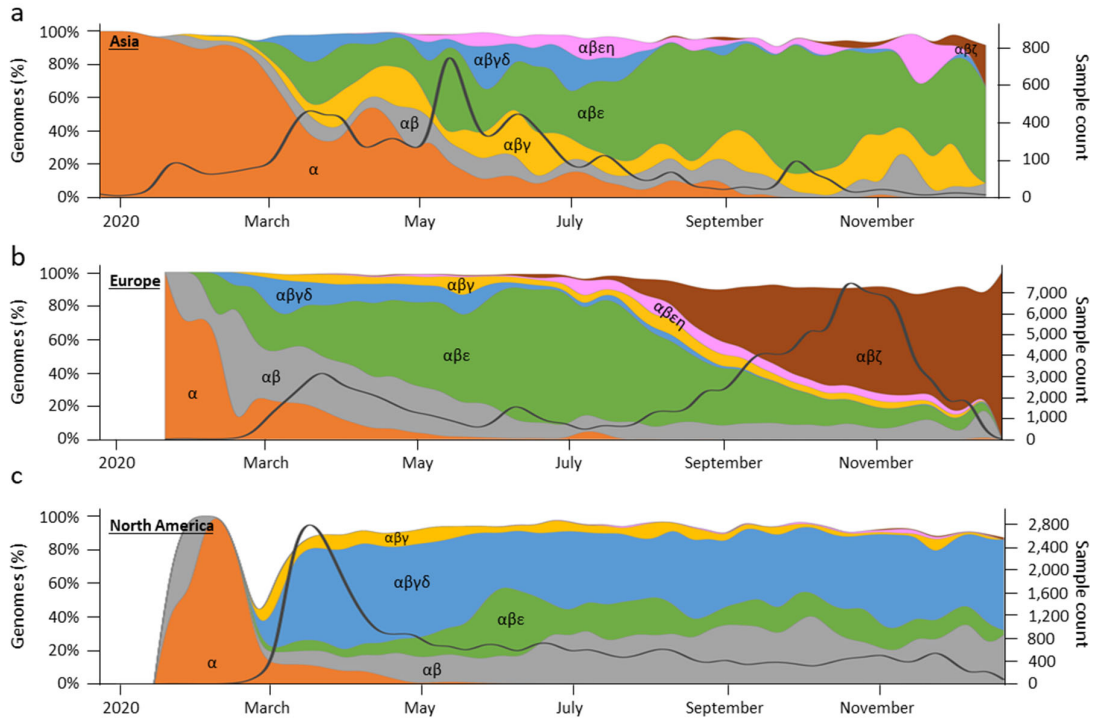


555

556

557 **Fig. 3. The backbone of SARS-CoV-2 mutational history.** The mutational history inferred was from (a) 29KG  
558 and (b) 68KG datasets. Major variants and their mutational pathways are shown in black, and minor variants and  
559 their mutational pathways are gray. Circle color marks the region where variants were sampled first. The 68KG  
560 dataset contains 12 additional variants and more than two times the genomes than the 29KG dataset.

561



562

563

564 **Fig. 4. Spatiotemporal dynamics of 172,480 SARS-CoV-2 genomes (December 2019-2020).** Spatiotemporal  
565 patterns of genomes mapped to lineages containing different combinations of major variants in (a) Asia, (b)  
566 Europe, and (c) North America. The number of genomes mapped to major variant lineages contains all of its offshoots, e.g.,  
567  $\alpha$  lineage contains all the genomes with  $\alpha_1 - \alpha_3$ ,  $\alpha_{1a} - \alpha_{1d}$ , and  $\alpha_{3a} - \alpha_{3j}$  variants only. The stacked graph area is  
568 the proportion of genomes mapped to the corresponding lineage. The solid black line shows the count of total  
569 genome samples. Spatiotemporal patterns in cities, countries, and other regions are available online at  
570 <http://sars2evo.datamonkey.org/>.

571

572

573 **Extended Data Table 1.** SARS-CoV-2 variants in 29KG dataset.

Mutant (major)	Mutant (minor)	Gene	Genomic Position	Nucleotide change	Amino acid change	Time (days)	Variants Frequency	Genomes mapped	First location
$\mu_1$		ORF1ab	2416	U>C		0	98.1%	0	China, Asia
$\mu_2$		ORF1ab	19524	U>C		0	98.6%	0	China, Asia
$\mu_3$		S	23929	U>C		0	98.4%	18	China, Asia
$\alpha_1$		ORF1ab	18060	U>C		0	95.1%	849	China, Asia
	$\alpha_{1a}$	N	28657	C>U		63	1.3%	2	France, Europe
	$\alpha_{1b}$	ORF1ab	9477	U>A	F>Y	63	1.2%	3	France, Europe
	$\alpha_{1c}$	N	28863	C>U	S>L	63	1.2%	5	France, Europe
	$\alpha_{1d}$	ORF3a	25979	G>U	G>V	63	1.2%	344	France, Europe
$\alpha_2$		ORF1ab	8782	U>C		0	91.0%	47	China, Asia
$\alpha_3$		ORF8	28144	C>U	S>L	0	90.8%	1115	China, Asia
	$\alpha_{3a}$	ORF1ab	1606	U>C		43	1.7%	501	United Kingdom, Europe
	$\alpha_{3b}$	ORF1ab	11083	G>U	L>F	24	9.2%	376	China, Asia
	$\alpha_{3c}$	N	28311	C>U	P>L	64	1.9%	3	South Korea, Asia
	$\alpha_{3d}$	ORF1ab	13730	C>U	A>V	71	1.8%	3	Taiwan/Malaysia, Asia
	$\alpha_{3e}$	ORF1ab	6312	C>A	T>K	71	1.7%	483	Taiwan/Malaysia, Asia
	$\alpha_{3f}$	ORF3a	26144	G>U	G>V	28	5.1%	121	China, Asia
	$\alpha_{3g}$	ORF1ab	14805	C>U		54	6.0%	334	United Kingdom, Europe
	$\alpha_{3h}$	ORF1ab	17247	U>C		64	2.0%	580	Switzerland, Europe
	$\alpha_{3i}$	ORF1ab	2558	C>U	P>S	54	1.7%	26	United Kingdom, Europe
	$\alpha_{3j}$	ORF1ab	2480	A>G	I>V	54	1.6%	462	United Kingdom, Europe
$\beta_1$		ORF1ab	3037	C>U		31	77.0%	11	China, Asia
$\beta_2$		S	23403	A>G	D>G	31	77.1%	36	China, Asia
$\beta_3$		ORF1ab	14408	C>U	P>L	41	76.9%	3032	Saudi Arabia, Middle East
	$\beta_{3a}$	ORF1ab	20268	A>G		64	5.7%	1213	Italy, Europe
	$\beta_{3b}$	N	28854	C>U	S>L	29	3.1%	527	China, Asia
	$\beta_{3c}$	ORF1ab	15324	C>U		29	2.3%	678	China, Asia
	$\beta_{3d}$	ORF3a	25429	G>U	V>L	77	1.7%	485	United Kingdom, Europe
	$\beta_{3e}$	N	28836	C>U	S>L	74	1.6%	3	Switzerland, Europe
	$\beta_{3f}$	ORF1ab	13862	C>U	T>I	74	1.6%	50	Switzerland, Europe
	$\beta_{3g}$	ORF1ab	10798	C>A	D>E	86	1.4%	414	United Kingdom, Europe
$\gamma_1$		ORF3a	25563	G>U	Q>H	41	29.8%	884	Saudi Arabia, Middle East
	$\gamma_{1a}$	ORF1ab	18877	C>U		41	4.0%	757	Saudi Arabia, Middle East
	$\gamma_{1b}$	M	26735	C>U		41	1.5%	439	Saudi Arabia, Middle East
$\delta_1$		ORF1ab	1059	C>U	T>I	54	23.0%	5157	Singapore, Asia
	$\delta_{1a}$	S	24368	G>U	D>Y	75	1.3%	389	Sweden, Europe
	$\delta_{1b}$	ORF8	27964	C>U	S>L	76	2.7%	790	USA, North America
	$\delta_{1c}$	ORF1ab	11916	C>U	S>L	72	1.6%	166	USA, North America
	$\delta_{1d}$	ORF1ab	18998	C>U	A>V	72	1.0%	305	USA, North America
$\epsilon_1$		N	28881	G>A	R>K	54	25.7%	2	United Kingdom, Europe
$\epsilon_2$		N	28882	G>A	R>K	54	25.7%	2	United Kingdom, Europe
$\epsilon_3$		N	28883	G>C	G>R	54	25.7%	5365	United Kingdom, Europe
	$\epsilon_{3a}$	ORF1ab	313	C>U		66	2.1%	608	USA, North America
	$\epsilon_{3b}$	ORF1ab	19839	U>C		64	1.5%	452	Switzerland, Europe
	$\epsilon_{3c}$	M	27046	C>U	T>M	69	1.6%	453	Worldwide
	$\epsilon_{3d}$	ORF1ab	10097	G>A	G>S	69	2.5%	5	Denmark, Europe
	$\epsilon_{3e}$	S	23731	C>U		69	2.5%	403	Denmark, Europe
	$\epsilon_{3f}$	N	28580	G>U	D>Y	69	1.2%	353	Chile, South America
$\nu_1$		ORF1ab	17858	A>G	Y>C	59	4.7%	32	USA, North America
$\nu_2$		ORF1ab	17747	C>U	P>L	59	4.7%	1374	USA, North America

Note.- Genomic locations correspond to those of the NCBI genome (GenBank ID: NC\_04551.2). Amino acid changes are shown for nonsynonymous variants.

574  
575  
576  
577

578 **Extended Data Table 2.** SARS-CoV-2 variants in the 68KG dataset.

Mutant (major)	Mutant (minor)	Gene	Genomic Position	Nucleotide change	Amino acid change	Time (days)	Variant Frequency	Genomes mapped	First location
$\mu_1$		ORF1ab	2416	U>C		0	98.4%	0	China, Asia
$\mu_2$		ORF1ab	19524	U>C		0	99.0%	18	China, Asia
$\mu_3$		S	23929	U>C		0	98.9%	0	China, Asia
$\mu_4$		ORF1ab	15933	U>C		0	98.8%	0	China, Asia
$\mu_5$		ORF8	27944	U>C		0	97.0%	0	China, Asia
$\mu_6$		ORF1ab	6286	U>C		0	95.6%	0	China, Asia
$\mu_7$		S	22444	U>C		0	98.7%	0	China, Asia
$\alpha_1$		ORF1ab	18060	U>C		0	97.3%	1114	China, Asia
	$\alpha_{1a}$	N	28657	C>U		63	1.0%	3	France, Europe
	$\alpha_{1b}$	ORF1ab	9477	U>A	F>Y	63	0.7%	3	France, Europe
	$\alpha_{1c}$	N	28863	C>U	S>L	63	0.7%	7	France, Europe
	$\alpha_{1d}$	ORF3a	25979	G>U	G>V	63	0.7%	451	France, Europe
$\alpha_2$		ORF1ab	8782	U>C		0	94.9%	51	China, Asia
$\alpha_3$		ORF8	28144	C>U	S>L	0	94.9%	1281	China, Asia
	$\alpha_{3a}$	ORF1ab	1606	U>C		43	0.9%	578	United Kingdom, Europe
	$\alpha_{3b}$	ORF1ab	11083	G>U	L>F	24	7.5%	417	China, Asia
	$\alpha_{3c}$	N	28311	C>U	P>L	64	1.4%	4	South Korea, Asia
	$\alpha_{3d}$	ORF1ab	13730	C>U	A>V	33	1.4%	5	China, Asia
	$\alpha_{3e}$	ORF1ab	6312	C>A	T>K	71	1.2%	767	Taiwan, Asia
	$\alpha_{3f}$	ORF3a	26144	G>U	G>V	28	3.0%	160	China, Asia
	$\alpha_{3g}$	ORF1ab	14805	C>U		54	3.7%	511	United Kingdom, Europe
	$\alpha_{3h}$	ORF1ab	17247	U>C		64	1.0%	682	Switzerland, Europe
	$\alpha_{3i}$	ORF1ab	2558	C>U	P>S	54	1.0%	44	United Kingdom, Europe
	$\alpha_{3j}$	ORF1ab	2480	A>G	I>V	54	1.0%	648	United Kingdom, Europe
$\beta_1$		ORF1ab	3037	C>U		31	87.2%	45	China, Asia
$\beta_2$		S	23403	A>G	D>G	31	87.2%	15	China, Asia
$\beta_3$		ORF1ab	14408	C>U	P>L	41	87.1%	4450	Saudi Arabia, Middle East
	$\beta_{3a}$	ORF1ab	20268	A>G		64	6.0%	2388	Italy, Europe
	$\beta_{3b}$	N	28854	C>U	S>L	29	4.5%	1782	China, Asia
	$\beta_{3c}$	ORF1ab	15324	C>U		29	2.2%	1463	China, Asia
	$\beta_{3d}$	ORF3a	25429	G>U	V>L	77	1.1%	719	United Kingdom, Europe
	$\beta_{3e}$	N	28836	C>U	S>L	74	0.8%	3	Switzerland, Europe
	$\beta_{3f}$	ORF1ab	13862	C>U	T>I	74	0.8%	85	Switzerland, Europe
	$\beta_{3g}$	ORF1ab	10798	C>A		86	0.6%	435	United Kingdom, Europe
$\gamma_1$		ORF3a	25563	G>U	Q>H	41	24.4%	1671	Saudi Arabia, Middle East
	$\gamma_{1a}$	ORF1ab	18877	C>U		41	4.2%	1201	Saudi Arabia, Middle East
	$\gamma_{1b}$	M	26735	C>U		41	2.7%	1784	Saudi Arabia, Middle East
$\delta_1$		ORF1ab	1059	C>U	T>I	54	17.6%	8284	Singapore, Asia
	$\delta_{1a}$	S	24368	G>U	D>Y	75	0.7%	466	Sweden, Europe
	$\delta_{1b}$	ORF8	27964	C>U	S>L	76	2.9%	1152	USA, North America
	$\delta_{1c}$	ORF1ab	11916	C>U	S>L	72	1.9%	807	USA, North America
	$\delta_{1d}$	ORF1ab	18998	C>U	A>V	72	0.7%	458	USA, North America
	$\delta_{1e}$	ORF1ab	10319	C>U	L>F	76	1.2%	799	USA, North America
$\zeta_1$		ORF1ab	445	U>C		179	4.4%	18	Netherlands, Europe
$\zeta_2$		M	26801	C>G		82	4.3%	7	Canada, North America
$\zeta_3$		S	22227	C>U	A>V	84	4.5%	1	Spain, Europe
$\zeta_4$		N	28932	C>U	A>V	96	4.4%	5	Portugal, Europe
$\zeta_5$		ORF10	29645	G>U	V>L	78	4.4%	2	Denmark, Europe
$\zeta_6$		ORF1ab	21255	G>C		80	4.4%	1557	USA, North America
$\zeta_7$		S	21614	C>U	L>F	79	2.5%	1442	United Kingdom, Europe

579

580

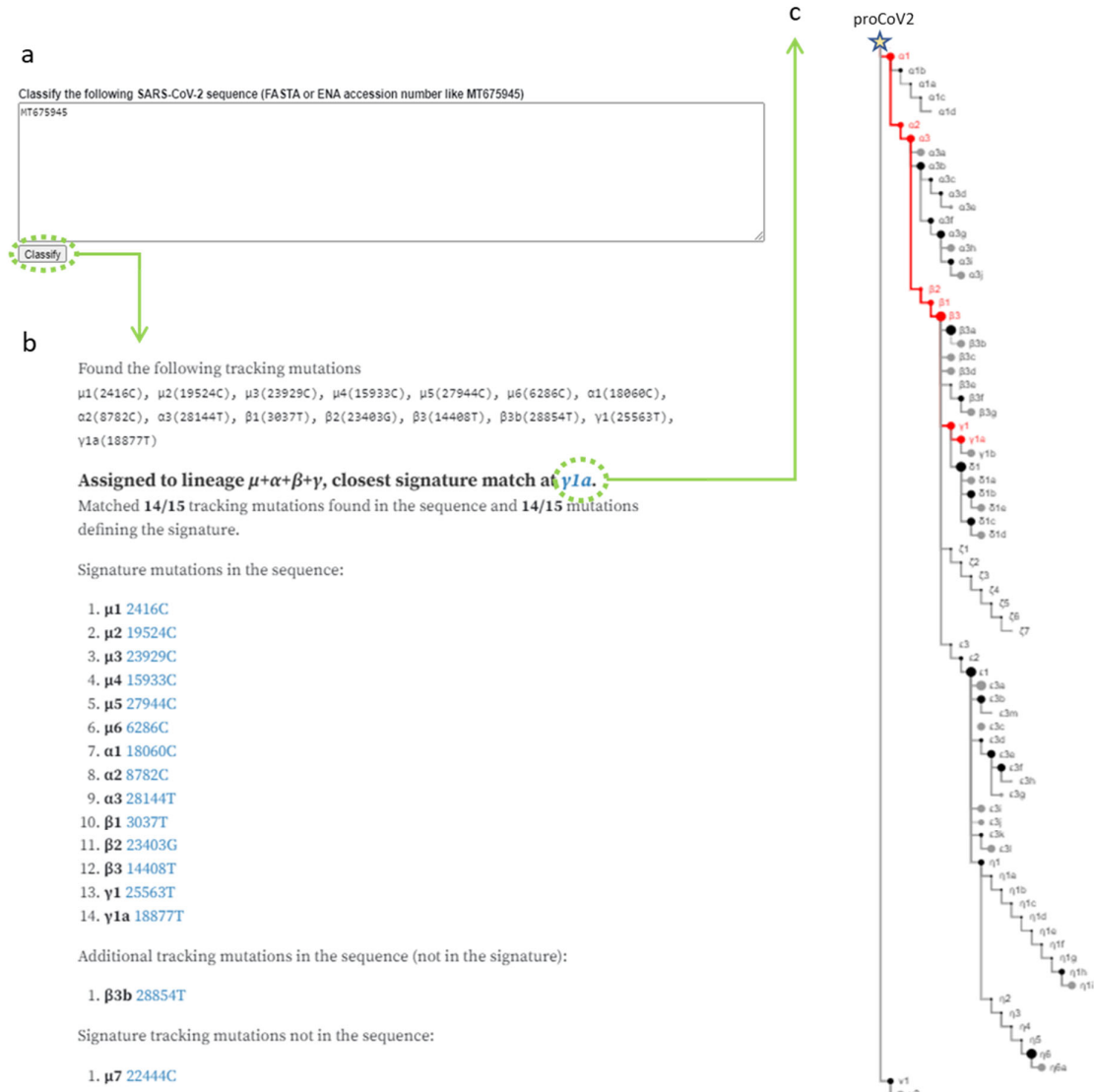
581 **Extended Data Table 2. SARS-CoV-2 variants in the 68KG dataset (continued).**

Mutant (major)	Mutant (minor)	Gene	Genomic Position	Nucleotide change	Amino acid change	Time (days)	Variant Frequency	Genomes mapped
ε <sub>1</sub>		N	28881	G>A	R>K	54	41.7%	5
ε <sub>2</sub>		N	28882	G>A	R>K	54	41.6%	0
ε <sub>3</sub>		N	28883	G>C	G>R	54	41.6%	13394
	ε <sub>3a</sub>	ORF1ab	313	C>U		64	2.4%	1630
	ε <sub>3b</sub>	ORF1ab	19839	U>C		64	2.9%	1227
	ε <sub>3c</sub>	M	27046	C>U	T>M	69	0.8%	548
	ε <sub>3d</sub>	ORF1ab	10097	G>A	G>S	69	3.2%	11
	ε <sub>3e</sub>	S	23731	C>U		69	3.2%	425
	ε <sub>3f</sub>	N	28580	G>U	D>Y	69	1.0%	678
	ε <sub>3g</sub>	ORF1ab	13536	C>U		69	1.6%	23
	ε <sub>3h</sub>	ORF1ab	4002	C>U	T>I	69	1.6%	1066
	ε <sub>3i</sub>	ORF1ab	10265	G>A	G>S	63	1.4%	879
	ε <sub>3j</sub>	S	21575	C>U	L>F	54	1.0%	248
	ε <sub>3k</sub>	S	21637	C>U		111	1.3%	873
	ε <sub>3l</sub>	ORF8	28169	A>G		103	1.3%	0
	ε <sub>3m</sub>	ORF1ab	16968	G>U		114	1.0%	702
η <sub>1</sub>		ORF1ab	1163	A>U	I>F	86	9.6%	339
	η <sub>1a</sub>	ORF1ab	14202	G>U		159	1.1%	7
	η <sub>1b</sub>	ORF1ab	19542	G>U	M>I	81	1.2%	23
	η <sub>1c</sub>	S	22388	C>U		90	1.2%	21
	η <sub>1d</sub>	N	29466	C>U	A>V	91	1.2%	4
	η <sub>1e</sub>	ORF1ab	19718	C>U	T>I	73	1.5%	23
	η <sub>1f</sub>	ORF3a	26060	C>U	T>I	92	1.2%	7
	η <sub>1g</sub>	N	29227	G>U		55	1.2%	24
	η <sub>1h</sub>	ORF1ab	3256	U>C		167	1.1%	0
	η <sub>1i</sub>	ORF1ab	5622	C>U	P>L	67	1.2%	775
η <sub>2</sub>		ORF1ab	18555	C>U		51	8.0%	25
η <sub>3</sub>		ORF1ab	16647	G>U		84	8.0%	8
η <sub>4</sub>		ORF1ab	7540	U>C		86	7.9%	0
η <sub>5</sub>		S	23401	G>A		86	7.9%	1
η <sub>6</sub>		S	22992	G>A	S>N	86	8.5%	4583
	η <sub>6a</sub>	S	22480	C>U		66	1.3%	878
v <sub>1</sub>		ORF1ab	17858	A>G	Y>C	59	2.6%	61
v <sub>2</sub>		ORF1ab	17747	C>U	P>L	59	2.5%	1677

582

583 Note.- Genomic locations correspond to those of the NCBI genome (GenBank ID: NC\_04551.2). Amino  
584 acid changes are shown for nonsynonymous variants.





594

595

596 **Extended Data Figure 2.** An example of sequence classification (ENA Accession MT675945) based on the  
 597 84 signature mutations (<http://sars2evo.datamonkey.org/>; "Classify your Sequence" option). (a) Input  
 598 window to provide identifiers of sequences to be classified (e.g., MT675945). (b) The input sequence is  
 599 classified into a mutational fingerprint. A list of mutations that are appeared in the input sequence is  
 600 shown in the output window. (c) A waterfall phylogeny shows the input sequence's location in the  
 601 phylogeny, which appears after clicking the closet signature matched mutation in panel b.