

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17

Ongoing Global and Regional Adaptive Evolution of SARS-CoV-2

Nash D. Rochman^{1,*}, Yuri I. Wolf¹, Guilhem Faure², Feng Zhang^{2,3,4,5,6} and Eugene V. Koonin^{1,*}

¹National Center for Biotechnology Information, National Library of Medicine, Bethesda, MD 20894

²Broad Institute of MIT and Harvard, Cambridge, MA 02142; ³Howard Hughes Medical Institute, Massachusetts Institute of Technology, Cambridge, MA 02139; ⁴McGovern Institute for Brain Research, Massachusetts Institute of Technology, Cambridge, MA 02139; ⁵Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, MA 02139; and ⁶Department of Biological Engineering, Massachusetts Institute of Technology, Cambridge, MA 02139

For correspondence: nash.rochman@nih.gov, koonin@ncbi.nlm.nih.gov

Keywords: SARS-Cov-2, phylogeny, ancestral reconstruction, epistasis, globalization

18 **Abstract**

19

20 Unprecedented sequencing efforts have, as of October 2020, produced nearly 200,000
21 genomes of severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), the virus
22 responsible for COVID-19. Understanding the trends in SARS-CoV-2 evolution is
23 paramount to control the pandemic, but analysis of this enormous dataset is a major
24 challenge. We show that the ongoing evolution of SARS-CoV-2 over the course of the
25 pandemic is characterized primarily by purifying selection but a small set of sites,
26 including spike 614 and nucleocapsid 203-204 appear to evolve under positive
27 selection. In addition to the substitutions in the spike protein, multiple substitutions in the
28 nucleocapsid protein appear to be important for SARS-CoV-2 adaptation to the human
29 host. The positively selected mutations form a strongly connected network of apparent
30 epistatic interactions and are signatures of major partitions in the SARS-CoV-2
31 phylogeny. These partitions show distinct spatial and temporal dynamics, with both
32 globalization and diversification trends being apparent.

33

34

35 **Main**

36

37 High mutation rates of RNA viruses(1) enable virus adaptation at a staggering pace.
38 Nevertheless, robust sequence conservation indicates that purifying selection is the
39 principal force shaping the evolution of virus populations(2,3,4,5). The fate of a novel
40 zoonotic virus is in part determined by the race between public health intervention and
41 viral diversification. Even intermittent periods of positive selection can result in lasting
42 immune evasion, leading to oscillations in the size of the susceptible population and
43 ultimately a regular pattern of repeating epidemics, as has been demonstrated for
44 Influenza(6,7,8).

45

46 During the current coronavirus pandemic, understanding the degree and dynamics of
47 the diversification of severe acute respiratory syndrome coronavirus 2 (SARS-Cov-2) is
48 essential for establishing a practicable, proportionate public health response, from
49 guidelines on isolation and quarantine to vaccination(9). To investigate evolution of
50 SARS-CoV-2, we aggregated all available SARS-Cov-2 genomes as of October 11,
51 2020, from the three principal repositories: Genbank(10), Gisaid(11), and CNCB(12).
52 From the total of 197,453 submissions in these databases, 82,592 unique SARS-Cov-2
53 genome sequences were identified, and 39,695 high quality sequences were
54 incorporated into a global multisequence alignment (MSA) consisting of the
55 concatenated open reading frames with stop codons trimmed. The vast majority of the
56 sequences excluded from the MSA were removed due to a preponderance of
57 ambiguous characters (see Methods). The sequences in the final MSA correspond to
58 73,236 isolates with date and location metadata.

59

60 Several methods for coronavirus phylogenetic tree inference have been tested(13,14).
61 The construction of a single high-quality tree from 40,000 30 kilobase (kb) sequences
62 using any of the existing advanced methods is computationally prohibitive. Therefore,
63 building on the available techniques, we assembled a set of maximally diverse subtrees
64 over a reduced alignment. This alignment was constructed mainly through the removal
65 of sites invariant up to the exception of a single sequence (see Methods). The omission

66 of these sites created additional sequence redundancy which was then eliminated by
67 removing additional sequences. The subtrees generated from the resulting MSA were
68 then used to constrain a single composite tree, both steps requiring a sophisticated
69 phylogeny reconstruction approach (see Methods). This composite tree reflects the
70 correct topology but has incorrect branch lengths and was in turn used to constrain a
71 global tree over the entire MSA (Fig. 1A). A complete reconstruction of ancestral
72 sequences was then performed by leveraging Fitch traceback(15) (see Methods),
73 enabling comprehensive identification of nucleotide and amino acid replacements
74 across the tree.

75
76 We identified 6 principal partitions within this tree, in general agreement with other
77 work(16, 17, 18), along with two region-specific clades within partitions 3 and 6 (Fig. 1A)
78 that, as discussed below, are important for the interpretation of the metadata. Given the
79 short evolutionary distances between SARS-CoV-2 isolates, the topology of the tree is a
80 cause of legitimate concern(14, 19, 20, 21). For the analyses presented below, we rely on
81 a single, explicit tree topology which is probably one of many equally likely
82 estimates(14). Therefore, we sought to validate the robustness of the major partitions of
83 the virus genomes using a phylogeny-free approach. To this end, pairwise Hamming
84 distances were computed for all sequences in the MSA and the resulting distance
85 matrix was embedded within a 3-dimensional subspace using classical multidimensional
86 scaling (Fig. 1B). In this embedding, all 6 partitions are nearly completely separated,
87 and the optimal clustering, determined by *k*-means, returned 4 categories (see
88 Methods, Fig. S1), grouping together partitions 1 and 2 as well as 4 and 5. These
89 findings indicate that an alternative tree with a comparable likelihood but a dramatically
90 different coarse-grain topology, most likely, cannot be constructed from this MSA.

91
92 Each of the 6 partitions can be characterized by a specific non-synonymous substitution
93 signature (Figs. 1C, S2), generally, corresponding to the most prominent non-
94 synonymous substitutions across the tree (Table S1), some of which are shared by two
95 or more partitions and appear independently many times, consistent with other
96 reports(22). The well known D614G substitution in the spike protein is part of these

97 signatures, and so are substitutions in two adjacent sites in the nucleocapsid protein
98 (see below). The rest of the signature sites are in the nonstructural proteins 1ab, 3a,
99 and 8 (Fig. 1C). The identification of these prevailing non-synonymous substitutions and
100 an additional set of frequent synonymous substitutions suggested that certain sites in
101 the SARS-CoV-2 genome might be evolving under positive selection. However,
102 uncovering the selective pressures affecting virus evolution was complicated by non-
103 negligible mutational biases. The distributions of the numbers of both synonymous and
104 non-synonymous substitutions across the genome were found to be substantially
105 overdispersed compared to both the Poisson and normal expectations (Fig. S3).
106 Examination of the relative frequencies of all 12 possible nucleotide substitutions
107 indicated a significant genome-wide excess of C to U mutations, approximately 3 fold
108 higher than any other nucleotide substitution, with the exception of G to U, as well as
109 some region-specific trends (Fig. S4).

110
111 Motivated by this observation, we compared the trinucleotide contexts of synonymous
112 and non-synonymous substitutions as well as the contexts of low and high frequency
113 substitutions. The context of high-frequency events, both synonymous and non-
114 synonymous, was found to be dramatically different from the background frequencies.
115 The NCN context (that is, all C->D mutations) harbors substantially more events than
116 other contexts (all 16 NCN triplets are within the top 20 most high-frequency-biased
117 ones, see Methods and Fig. S5) and is enriched in mutations uniformly across the
118 genome including both synonymous and non-synonymous sites as well as low and high
119 frequency sites. This pattern suggests a mechanistic bias of the coronavirus RNA-
120 dependent RNA polymerase (RdRP). Evidently, such a bias that increases the
121 likelihood of observing multiple, independent mutations in the NCN context complicates
122 the detection of selection pressures. However, whereas all the sites with an excess of
123 synonymous nucleotide substitutions are NCN and thus can be inferred to originate
124 from the mutational bias, this is not the case for non-synonymous substitutions,
125 suggesting that at least some of these are driven by other mechanisms. Thus, we
126 excluded all synonymous substitutions and the non-synonymous substitutions with the

127 NCN context from further consideration as candidate sites evolving under positive
128 selection.

129
130 Beyond this specific context, the presence of any hypervariable sites complicates the
131 computation of the dN/dS ratio, the gauge of protein-level selection(23), which requires
132 enumerating the number of synonymous and non-synonymous substitutions within each
133 gene. Hypervariable sites bias this analysis, and therefore, we used two methods to
134 ensure reliable estimation of dN/dS . For each protein-coding gene of SARS-CoV-2
135 (splitting the long orf1ab into 15 constituent non-structural proteins), we obtained both a
136 maximum likelihood estimate of dN/dS across 10 sub-alignments and an approximation
137 computed from the global ancestral reconstruction (see Methods). This approach was
138 required due to the size of the alignment, over which a global maximum likelihood
139 estimation would be computationally prohibitive. Despite considerable variability among
140 the genes, we obtained estimates of substantial purifying selection ($0.1 < dN/dS < 0.5$)
141 across the majority of the genome(Fig. S6), with a reasonable agreement between the
142 two methods. This estimate is compatible with previous demonstrations of purifying
143 selection among diverse RNA viruses(3) affecting about 50% of the sites surveyed or
144 more(2).

145
146 Thus, the evolution of SARS-CoV-2 appears to be primarily driven by substantial
147 purifying selection. However, a small ensemble of non-synonymous substitutions
148 appeared to have emerged multiple times, independently, and were not subject to an
149 overt mechanistic bias. Due to the existence of many equally likely trees, in principle, in
150 one or more of such trees, any of these mutations could resolve to a single event.
151 However, such a resolution would be at the cost of inducing multiple parallel
152 substitutions for other mutations, and thus, we conclude that more than 100 codons in
153 the genome have undergone multiple parallel mutations in the course of SARS-CoV-2
154 evolution.

155
156 One immediate explanation of this observation is that these sites evolve under positive
157 selection. The possible alternatives could be that these sites are mutational hotspots or

158 that the appearance of multiple parallel mutations was caused by numerous
159 recombination events (either real or artifacts caused by incorrect genome assembly
160 from mixed infections) in the respective genomic regions. Contrary to what one would
161 expect under the hotspot scenario, we found that codons harboring many synonymous
162 substitutions tend to harbor few non-synonymous substitutions, and vice versa (Fig. S7
163 A). Although when a moving average with increasing window size was computed, this
164 relationship reversed (Fig. S7 B&C), the correlation between synonymous and non-
165 synonymous substitutions was weak. Most sites in the virus genome are highly
166 conserved, those sites that harbor the highest number of mutations tend to reside in
167 conserved neighborhoods, and the local fraction of sites that harbor at least one
168 mutation strongly correlates with the moving average (Fig. S8). Thus, whereas our
169 observations indicate that SARS-CoV-2 genomes are subject to diverse site-specific
170 and regional selection pressures, we did not detect regions of substantially elevated
171 mutation or recombination.

172
173 Given the widespread purifying selection, substantially relaxed selection at any site is
174 expected to permit multiple, parallel non-synonymous mutations to the same degree
175 that any site harbors multiple, parallel synonymous mutations. Thus, seeking to identify
176 sites subject to positive selection, we focused only on those non-synonymous
177 substitutions that independently occurred more frequently than 95% of all synonymous
178 substitutions excluding the mutagenic NCN context (see Methods). Most if not all sites
179 in the SARS-CoV-2 genome that we found to harbor such frequent, parallel non-
180 synonymous substitutions not subject to the restrictions discussed above are likely to
181 evolve under positive selection (Fig. 1D, Table S2).

182
183 Having identified the set of potential positively selected residues, we examined the tree
184 for evidence of epistasis⁽²⁴⁾ (see Methods) among these sites and revealed a network
185 of co-occurring substitutions suggestive of epistatic interactions (Fig. 1E, Table S3).
186 Strikingly, both D614G in the spike (S) protein and two adjacent substitutions in the
187 nucleocapsid (N) protein, R203K and G204R, are associated with exceptionally many
188 interactions, forming the two main hubs of the network. Spike D614G appears to boost

189 the infectivity of the virus, possibly, by increasing the binding affinity between the spike
190 protein and the cell surface receptor of SARS-CoV-2, ACE2 (25). The non-synonymous
191 mutations S|L54H,Q677H, A879S, and V1176F that are strongly linked to D614G in the
192 network are in the spike protein itself although examination of the spike structure does
193 not reveal direct physical interactions among these residues (Fig. S9). Another
194 substitution, S|H49Y, with a weaker statistical association to S|D614G, involves a site
195 that is close to 614 in the structure (Fig. S9). Conceivably, by increasing the receptor
196 affinity, the D614G substitution in the spike protein opens up new adaptive routes for
197 later steps in the viral lifecycle, but the specific mechanisms remain to be investigated
198 experimentally.

199
200 Of further interest were four substitutions within the spike protein that have been
201 observed in mink populations: H69del/V70del, Y453F, I692V and M1229I(26-27). I692V
202 never occurred within our tree, and this codon harbors many synonymous substitutions
203 atc(I) to att(I), suggestive of a mutational hotspot. In contrast, 69-, 70-, 453F, and 1229I
204 all appeared multiple times independently throughout the tree. These substitutions were
205 widely spread over the tree although all occurred close to tree tips and none passed our
206 criteria for positive selection. Two sites within the receptor-binding domain (RBD), N331
207 and N343, have been shown to be important for the maintenance of infectivity(28). As
208 could be expected, these amino acid residues are invariant. Four more substitutions in
209 the RBD, among others, N234Q, L452R, A475V, and V483A, have been demonstrated
210 to confer antibody resistance(28). We found N234 to be invariant, whereas a few L452R
211 and A475V mutations close to tree tips were detected as well as two independent
212 V483A mutations at the base of slightly larger clades within partition 1. None of these
213 sites passed our criteria for positive selection.

214
215 The two adjacent amino acid replacements in the N protein, R(agg)203K(aaa) and
216 G(gga)204R(cga), appear together 9 times. Both sites are likely to evolve under positive
217 selection and are adjacent to a third such site, S(agt)202N(aat). The replacements
218 R(agg)203K(aaa) and G(gga)204R(cga) occur via three adjacent nucleotide
219 substitutions which strongly suggests a single mutational event. Evolution of beta-

220 coronaviruses with high case fatality rates including SARS-CoV-2 was accompanied by
221 accumulation of positive charges in the N protein that are thought to enhance its
222 transport to the nucleus(29). Although positions 202-204 are outside the known nuclear
223 localization signals, the substitutions in these positions are statistically associated with
224 the Q229H substitution in the N protein which occurs in a site known to be responsible
225 for nuclear shuttling (30). Furthermore, the rapid rise of the A220V substitution in the N
226 protein (excluded from considered as a candidate for positive selection in our analysis
227 due to its NCN context) in a European cohort during the summer of 2020 might be
228 related to a transmission advantage of the variant harboring this substitution(31).
229 Conceivably, the substitutions in these two adjacent sites, in particular
230 G(gga)204R(cga), which increases the positive charge, might contribute to the nuclear
231 localization of the N protein as well. This highly unusual cluster of three putative
232 positively selected amino acid substitutions in the N protein is a strong candidate for
233 experimental study that could illuminate the evolution of SARS-CoV-2 pathogenicity.

234
235 Orf3a|Q57H is a third hub in the network and, although not considered a candidate for
236 positive selection in our analysis due to its NCN context, ORF8 S84L is a hub in the
237 larger epistatic network including all strongly associated residues (Fig. S10). Also of
238 interest in this larger network is S|N439K that is linked to S|D614G and, despite its NCN
239 context, is potentially subject to positive selection having been demonstrated to enable
240 immune escape(32).

241
242 In addition to the recurrent missense mutations that are likely to evolve under positive
243 selection, we identified numerous nonsense mutations (Table S4), the most frequent
244 one, orf8|18(UAA), appearing in 66 unique genome sequences. These nonsense
245 substitutions, apparently, resulting in truncated proteins, occur almost exclusively within
246 the minor SARS-CoV-2 ORFs. ORF8 has been implicated in the modulation of host
247 immunity by SARS-CoV-2, so these truncations might play a role in immune
248 evasion(33-34).

249

250 Epistasis in RNA virus evolution, as demonstrated for influenza, can constrain the
251 evolutionary landscape as well as promote compensatory variation in coupled sites,
252 providing an adaptive advantage which would otherwise confer a prohibitive fitness
253 cost(35). Because even sites subject to purifying selection(36) can play an adaptive role
254 through interactions with other residues in the epistatic network, the network presented
255 here (Fig. 1E) likely underrepresents the extent of epistatic interactions occurring during
256 SARS-CoV-2 evolution. The early evolutionary events that shaped the epistatic network
257 conceivably laid the foundation for the diversification of the virus relevant to virulence,
258 immune evasion, and transmission. Recent analysis of within-patient genetic diversity
259 has shown that the most common mutations are highly diverse within individuals(37).
260 Such diversity could either result from multiple infections, or otherwise, could present
261 evidence of an even greater role of positive selection affecting a larger number of sites
262 than inferred from our tree (Fig. 1D). Similarly to the case of Influenza, positive selection
263 on these sites could drive virus diversification and might support a regular pattern of
264 repeat epidemics with grave implications for public health. An analysis of the
265 relationships between the sequencing date and location of each isolate and its position
266 within the tree can determine whether diversification is already apparent within the
267 evolutionary history of SARS-CoV-2.

268
269 We first demonstrated a strong correlation between the sequencing date of SARS-CoV-
270 2 genomes and the distance to the tree root (Fig. S11), indicating a sufficiently low level
271 of noise in the data for subsequent analyses. Although examination of the global
272 distribution of each of the 6 major SARS-CoV-2 partitions (Figs. S12-13) indicates
273 considerable regional diversification, this variation could partly result from time-
274 dependent fluctuations (Fig. 2). The beginning of the pandemic is primarily
275 characterized by the global extinction of partitions 1 and 2 which were dominant initially,
276 through the beginning of February, in all four regions encompassing the vast majority of
277 the available sequences (United States, Europe, Asia, and Australia/New Zealand).
278 Notably, partitions 1 and 2 lack the substitution SJD614G which is the consensus for all
279 other partitions.

280

281 The four “late” partitions began rising to prominence thereafter, with different regional
282 trends. A common global trend is the increasing prevalence of partition 6, the only one
283 in which the adjacent substitutions N|R203K and G204R belong to the consensus. From
284 April to August, partition 6 grew in prominence in all four regions, the only partition to do
285 so (Figs. 2 and 3A, S14). This trend is most dramatic in Australia/New Zealand where
286 all other partitions apparently went extinct by the beginning of June. Notably, most
287 sequences from this region form a clade offset by a long branch within partition 6
288 (mutational signature shown in Fig. S15 notably including S|S477N). Although the rise
289 of partition 6 in the US has been slow and unsteady, perhaps, the starkest contrast to
290 the global ascension of partition 6 is the resurgence of partition 3 within Europe.
291 Partition 3 seems to be poised to again become dominant in Europe after partition 6 had
292 risen to greater than 80% frequency in July. In this case, as in the case of partition 6 in
293 Australia/New Zealand, clustering of sequences offset by a long branch within partition
294 3 is observed. The non-synonymous substitution signature of this clade (Fig. S16)
295 includes N|A220V such that nearly all European isolates belong to partition 6 and
296 harbor N|R203K&G204R or belong to partition 3 and harbor N|A220V by the beginning
297 of September. Motivated by this finding, we examined the mutational signature
298 characterizing isolates collected after July 17 from the US, the region with the lowest
299 fraction of partition 6 isolates, compared with those collected after July 17 from the
300 remaining three regions (Fig. S17). Focusing only on substitutions in the N protein (Fig.
301 S18), three additional mutations in the vicinity of motifs implicated in nuclear
302 localization, N|P199L,S194L,T205I, were identified to be prevalent within the US (Fig.
303 S19). Although only few isolates from Asia after late August are available within this
304 dataset, the rapid decrease in partition 6 frequency corresponds to a rapid increase in
305 N|S194L within the region. Taken together, isolates with at least 1 of these 5 mutations
306 compose the majority in all regions by the beginning of August (Fig. 2E), strongly
307 suggesting that this region of the N protein plays a prominent role in the adaptation of
308 SARS-CoV-2 to human hosts.

309
310 Despite these global trends, regional dynamics make it difficult to assess whether
311 partition 6 or, perhaps, a distinct clade within partition 3 are indeed more fit than the

312 other “late” partitions and whether SARS-CoV-2 is beginning to regionally diversify. To
313 better assess the extent of regional divergence, we constructed two diversity measures,
314 one partition-dependent and the other, partition-independent. First, we considered 3
315 groups of virus genomes: partitions 1&2 (without the consensus substitution S|D614G),
316 partitions 3,4,5, and partition 6 (with consensus substitutions N|R203K&G204R), and
317 computed the Hellinger distance of this three-group frequency distribution between all
318 pairs of regions over a sliding window as a function of time (Fig. 3B). Next, we sampled
319 pairs of isolates both within and between regions, computed the mean tree-distance
320 between pairs over the same sliding window (Figs. 3C, S16), and examined the ratio of
321 inter-regional and intra-regional tree-distances (see Methods).

322
323 Both measures reveal three principal stages of the pandemic. 1) The beginning of the
324 year, through early February, when partitions 1 and 2 were dominant, was marked by
325 regional diversification. 2) The extinction of the two “early” partitions signaled
326 globalization and the spread of the mutation S|D614G from February through the end of
327 May. By June, Australia/New Zealand began diverging from the rest of the globe and
328 the remaining regions began to modestly diverge from one another as well. These
329 trends are superimposed over steady, substantial intra-regional diversification (Fig.
330 S16).

331
332 Such diversification of the virus could potentially pose problems for both testing and
333 vaccine development. Substitutions in the E protein have already been demonstrated to
334 interfere with a common PCR assay(38). Generally, ORF1ab is more conserved than
335 the spike protein, which itself is more conserved than the remaining ORFs (Figs. S3-4).
336 Using our SARS-CoV-2 MSA, we surveyed 10 regions from ORF1ab(5), N(4), and E(1)
337 genes that are commonly used within PCR assays (39) for substitutions relative to the
338 reference sequence. Among the nearly 40,000 genome sequences, there were
339 hundreds to thousands of nucleotide substitutions in each of these regions but those in
340 ORF1ab were markedly less variable than those in N (Supplementary table 5) with one
341 region in N demonstrating variability in nearly half of all isolates. It can be expected that

342 most targets within the polyprotein will remain subject to the fewest polymorphism-
343 induced false negatives even as the virus continues to diversify.

344

345 Of the 9 vaccine candidates currently in phase 3 trials(40), three are inactivated whole-
346 virus (Sinovac, Wuhan Institute of Biological Products/Sinopharm, Beijing Institute of
347 Biological Products/Sinopharm); five utilize the entire spike protein as the antigen
348 (Moderna/NIAID, CanSino Biological Inc./Beijing Institute of Biotechnology, University of
349 Oxford/AstraZeneca, Gamaleya Research Institute, Janssen Pharmaceutical
350 Companies) and one utilizes only the RBD (Pfizer/Fosun Pharma/BioNTech). In addition
351 to the greater sequence conservation of the spike protein relative to all other ORFs
352 outside of the polyprotein, it is the principal host-interacting protein of SARS-CoV-2,
353 making both the whole protein and the RBD obvious antigenic candidates. Most
354 mutations in the RBD were demonstrated to decrease infectivity; however, some
355 conferred resistance to neutralizing antibodies(27). We only identified one mutation in
356 the RBD that was both observed in greater than 1% of the tree branches and passed
357 our criteria for positive selection, S|S477N. Should such mutations become more
358 prominent, the RBD might prove a less effective antigen than the whole spike protein.

359

360 To summarize, from these findings, it is clear that, despite dramatically reduced
361 travel(41), the evolution of SARS-Cov-2 is at least partly shaped by globalizing factors,
362 such as the increased virus fitness conferred by S|D614G, N|R203K&G204R, and other
363 positively selected substitutions. There is no strong evidence of “speciation”, that is,
364 formation of stable, diverging variants, a finding that bodes well for a successful
365 vaccination campaign in the midterm. Nevertheless, it is equally clear that SARS-CoV-2
366 has the capacity to diversify, posing a risk of virus escape from immunity and hence
367 repeat epidemics.

368

369 **Author contributions**

370 EVK initiated the project; NR and GF collected data; NR, GF, YIW, FZ and EVK
371 analyzed data; NR and EVK wrote the manuscript that was edited and approved by all
372 authors.

373 **Acknowledgements**

374 The authors thank Koonin group members for helpful discussions. NR, YIW
375 and EVK are supported by the Intramural Research Program of the
376 National Institutes of Health (National Library of Medicine).

377

378 **References**

379

380 [1] J.W. Drake, J.J. Holland. Mutation rates among RNA viruses. *Proceedings of the*
381 *National Academy of Sciences* **96.24**, 13910-13913 (1999).

382

383 [2] J.O. Wertheim, and S.L. Kosakovsky Pond. Purifying selection can obscure the
384 ancient age of viral lineages. *Molecular biology and evolution* **28.12**, 3355-3365 (2011)

385

386 [3] G.M. Jenkins et al. Rates of molecular evolution in RNA viruses: a quantitative
387 phylogenetic analysis. *Journal of molecular evolution* **54.2**, 156-165 (2002)

388

389 [4] E.C. Holmes. Patterns of intra-and interhost nonsynonymous variation reveal strong
390 purifying selection in dengue virus. *Journal of virology* **77.20**, 11296-11298 (2003)

391

392 [5] G. Jerzak, et al. Genetic variation in West Nile virus from naturally infected
393 mosquitoes and birds suggests quasispecies structure and strong purifying selection.
394 *The Journal of general virology* **86.Pt 8**, 2175 (2005)

395

396 [6] Y.I. Wolf, et al. Long intervals of stasis punctuated by bursts of positive selection in
397 the seasonal evolution of influenza A virus. *Biology direct* **1.1**, 34 (2006)

398

399 [7] R.M. Bush et al. Predicting the Evolution of Human Influenza A. *Science* **286.5446**
400 (1999).

401

402 [8] R.M. Bush et al. Positive selection on the H3 hemagglutinin gene of human influenza
403 virus A. *Molecular biology and evolution* **16.11** (1999).

404

405 [9] A. Koirala, et al. Vaccines for COVID-19: The current state of play. *Paediatric*
406 *respiratory reviews* **35**, 43-49 (2020)

407

- 408 [10] D.A. Benson, et al. GenBank. *Nucleic acids research* **41.D1**, D36-D42 (2012)
409
- 410 [11] S. Elbe, G. Buckland-Merrett. Data, disease and diplomacy: GISAID's innovative
411 contribution to global health. *Global Challenges* **1.1**, 33-46 (2017)
412
- 413 [12] W. Zhao et al. The 2019 novel coronavirus resource. *Hereditas* **42.2**, 212-221
414 (2020)
415
- 416 [13] R. Lanfear. *A global phylogeny of SARS-CoV-2 from GISAID data, including*
417 *sequences deposited up to 20-August-2020. Zenodo* (2020). DOI:
418 10.5281/zenodo.3958883
419
- 420 [14] B. Morel et al. Phylogenetic analysis of SARS-CoV-2 data is difficult. *bioRxiv*
421 (2020).
422
- 423 [15] W.M. Fitch. Toward defining the course of evolution: minimum change for a specific
424 tree topology. *Systematic Biology* **20.4**, 406-416 (1971)
425
- 426 [16] S. Kumar et al. An evolutionary portrait of the progenitor SARS-CoV-2 and its
427 dominant offshoots in COVID-19 pandemic. *bioRxiv* (2020).
428
- 429 [17] P. Forster et al. Phylogenetic network analysis of SARS-CoV-2 genomes.
430 *Proceedings of the National Academy of Sciences* **117.17**, 9241-9243 (2020)
431
- 432 [18] N.M. Fountain-Jones et al. Emerging phylogenetic structure of the SARS-CoV-2
433 pandemic. *bioRxiv* (2020).
434
- 435 [19] C. Mavian et al. Sampling bias and incorrect rooting make phylogenetic network
436 tracing of SARS-COV-2 infections unreliable. *Proceedings of the National Academy of*
437 *Sciences* **117.23**, 12522-12523 (2020)
438

- 439 [20] S.J. Sánchez-Pacheco et al. Median-joining network analysis of SARS-CoV-2
440 genomes is neither phylogenetic nor evolutionary. *Proceedings of the National*
441 *Academy of Sciences* **117.23**, 12518-12519 (2020)
442
- 443 [21] L. Pipes et al. Assessing uncertainty in the rooting of the SARS-CoV-2 phylogeny.
444 *bioRxiv* (2020).
445
- 446 [22] L. van Dorp et al. Emergence of genomic diversity and recurrent mutations in
447 SARS-CoV-2. *Infection, Genetics and Evolution* **104351** (2020)
448
- 449 [23] Z. Yang and N. Goldman. A codon-based model of nucleotide substitution for
450 protein-coding DNA sequences. *Molecular biology and evolution* **11.5** (1994).
451
- 452 [24] N.D. Rochman, Y.I. Wolf, E.V. Koonin. Deep phylogeny of cancer drivers and
453 compensatory mutations. *Communications Biology* **3.1**, 1-11 (2020)
454
- 455 [25] B. Korber et al. Tracking changes in SARS-CoV-2 Spike: evidence that D614G
456 increases infectivity of the COVID-19 virus. *Cell* **182.4**, 812-827 (2020)
457
- 458 [26] Statens Serum Institut. Mutations in the mink virus.
459 <https://www.ssi.dk/aktuelt/nyheder/2020/mutationer-i-minkvirus>
460
- 461 [27] B. Munnink et al. Transmission of SARS-CoV-2 on mink farms between humans
462 and mink and back to humans. *Science* **eabe5901** (2020)
463
- 464 [28] L. Qianqian et al. The Impact of Mutations in SARS-CoV-2 Spike on
465 Viral Infectivity and Antigenicity. *Cell* **182**, 1284–1294 (2020)
466
- 467 [29] A.B. Gussow et al. Genomic determinants of pathogenicity in SARS-CoV-2 and
468 other human coronaviruses. *Proceedings of the National Academy of Sciences* (2020).
469

- 470 [30] K.A. Timani et al. Nuclear/nucleolar localization properties of C-terminal
471 nucleocapsid protein of SARS coronavirus. *Virus research* **114.1-2**, 23-34 (2005)
472
- 473 [31] E. Hodcroft et al. Emergence and spread of a SARS-CoV-2 variant through Europe
474 in the Summer of 2020. *medRxiv* (2020).
475
- 476 [32] E. Thomson et al. The circulating SARS-CoV-2 spike variant N439K maintains
477 fitness while evading antibody-mediated immunity. *bioRxiv* (2020).
478
- 479 [33] Y. Zhang et al. The ORF8 Protein of SARS-CoV-2 Mediates Immune Evasion
480 through Potently Downregulating MHC-I. *bioRxiv* (2020).
481
- 482 [34] L. Zinzula Lost in deletion: The enigmatic ORF8 protein of SARS-CoV-2.
483 *Biochemical and Biophysical Research Communications* (2020).
484
- 485 [35] L.I. Gong, M.A. Suchard, J.D. Bloom. Stability-mediated epistasis constrains the
486 evolution of an influenza protein. *Elife* **2**, e00631 (2013)
487
- 488 [36] S. Kryazhimskiy et al. Prevalence of epistasis in the evolution of influenza A surface
489 proteins. *PLoS Genet* **7.2**, e1001301 (2011)
490
- 491 [37] J. Kuipers et al. Within-patient genetic diversity of SARS-CoV-2. *bioRxiv* (2020).
492
- 493 [38] M. Artesi et al. A recurrent mutation at position 26340 of SARS-CoV-2 is associated
494 with failure of the E gene quantitative reverse transcription-PCR utilized in a commercial
495 dual-target diagnostic assay. *Journal of clinical microbiology* **58.10** (2020).
496
- 497 [39] E. Ortiz-Prado et al. Clinical, molecular, and epidemiological characterization of the
498 SARS-CoV-2 virus and the Coronavirus Disease 2019 (COVID-19), a comprehensive
499 literature review. *Diagnostic Microbiology and Infectious Disease*. **115094** (2020).
500

501 [40] Y. Dong et al. A systematic review of SARS-CoV-2 vaccine candidates. *Signal*
502 *Transduction and Targeted Therapy* **5**, 237 (2020).

503

504 [41] S. Lai et al. Assessing the effect of global travel and contact reductions to mitigate
505 the COVID-19 pandemic and resurgence. *medRxiv* (2020).

506

507 **Methods**

508

509 **Multiple alignment of SARS-CoV-2 genomes**

510 All available SARS-CoV-2 genomes as of October 11, 2020 were retrieved from the
511 Genbank(10), Gisaid(11), and CNCB(12) datasets. Sequences with apparent anomalies
512 (sequence inversion etc.) were immediately discarded. Sequences were harmonized to
513 DNA (e.g. U was transformed to T to amend software compatibility) and clustered
514 according to 100% identity with no coverage threshold using CD-HIT(42-43), with
515 ambiguous characters masking. All characters excepting ACGT were considered
516 ambiguous. The least ambiguous sequence from each cluster was selected and
517 sequences shorter than 25120 nucleotides were discarded.

518 Exterior ambiguous characters (preceding/succeeding the first/last defined nucleotide)
519 were removed and sequences with more than 10 remaining, interior, ambiguous
520 characters were discarded. The remaining sequences were aligned using multi-
521 threaded MAFFT(44) with 220 cores (--thread 220) and 3.8Tb of RAM to maintain
522 usage of the normal DP algorithm (44) (--nomemsave). Sequences sourced from non-
523 human hosts were manually identified from the metadata and those excluded at the
524 previous step were added to the alignment using MAFFT, maintaining the number of
525 columns in the original alignment (specifying --keeplength), again on 220 cores.

526 Sites corresponding to protein-coding open reading frames were then mapped to the
527 alignment from the reference sequence NC_045512.2 excluding stop codons as follows:
528 266-13468+13468-21552, orf1ab; 21563-25381, S; 25393-26217, orf3a; 26245-26469,
529 E; 26523-27188, M; 27202-27384, orf6; 27394-27756, orf7a; 27756-27884, orf7b;
530 27894-28256, orf8; and 28274-29530, N. The remaining sites were discarded.

531 The resulting alignment contained out-of-frame gaps. Gaps in the reference sequence
532 were found to correspond to gaps in all but fewer than 1% of the remaining sequences.
533 These sites were discarded. The remaining gaps shorter than three nucleotides were
534 replaced with the ambiguous character, N. Longer gaps were shifted into frame and
535 padded with ambiguous characters on either end of the gap, minimizing the number of
536 sites altered.

537 A fast, approximate tree was then built using FastTree(45) (parameters: -nt -gtr -gamma
538 -nosupport -fastest) to unambiguously define two clusters of sequences: an outgroup
539 consisting of 14 sequences sourced from non-human hosts prior to 2020 and the main
540 group. The tree construction requires the resolution of very short branch lengths which
541 makes it necessary to compile FastTree at double precision. Outliers from the remaining
542 sequences were then identified based on the Hamming distance (excluding gaps and
543 ambiguous characters) to the nearest neighbor, the Hamming distance to the
544 consensus, and the degree to which those substitutions relative to consensus were
545 clustered in the genome. At this step, 21 sequences were removed.

546 The resulting alignment, consisting of 39,695 sequences and 29,119 sites, was
547 maintained for the construction of the global tree and ancestral sequence
548 reconstruction. In an effort to minimize the impact of sequencing error on the tree
549 topology, as well as to decrease computational costs, a reduced alignment was then
550 constructed through the removal of 1) invariant sites, 2) sites invariant with the
551 exception of a single sequence, and 3) sites invariant throughout the main group with
552 the exception of at most one sequence representing each minority nucleotide.
553 Removing these sites created substantial redundancy, so a representative sequence
554 was selected for each cluster of 100% identity to yield an alignment consisting of 34,685
555 sequences and 10,131 sites.

556

557 Tree Construction

558 We sought to optimize tree topology with IQ-TREE(46); however, we found building the
559 global tree to be computationally prohibitive, and thus, we proceeded to subsample the
560 main group alignment as follows. First, a core set of maximally diverse sequences is
561 selected. The set is initialized with a pair of sequences: a sequence maximizing the
562 number of substitutions relative to consensus and a paired sequence which maximizes
563 the Hamming distance to itself. Sequences are then added to this core set one at a time
564 maximizing the minimum Hamming distance to any representative of the set until N
565 sequences are incorporated. Next, $\text{ceil}(L/(M - N))$ resulting sets are initialized with this
566 core set where M is the target number of sequences and L is the total number of
567 sequences in the alignment (34,685). Then, sequences that have not yet been
568 incorporated into any resulting set are added to each resulting set, again one at a time,
569 maximizing the minimum distance to any representative of the set until M sequences
570 are incorporated. The order of the resulting sets is randomized at each iteration without
571 repeats. Once every (main group) sequence has been incorporated into at least one
572 resulting set, sequences are randomly incorporated into each set until every set
573 contains M sequences. Finally, the outgroup is added to each resulting set. We chose
574 $M=3,000$ in an effort to optimize computational efficiency and $N=300$. Note that while
575 increasing N increases the number of sets required for alignment coverage, and thus
576 compute time, insufficient overlap between the sequences assigned each sub-alignment

577 greatly affects the results of subsequent steps. We found $N=100$ to be too low to
578 effectively constrain the global tree for this dataset.

579 A tree was then built, using IQ-TREE, for each resulting set, with the evolutionary model
580 fixed to GTR+F+G4 and the minimum branch length decreased from the default $10e-6$
581 to $10e-7$, according to the results of previous parameter studies(14). These trees were
582 then converted into constraint files and merged to generate a single global constraint file
583 for use within FastTree (parameters: -nt -gtr -gamma -cat 4 -nosupport -constraints).

584 The remaining sequences excluded from this tree were then reintroduced as unresolved
585 multifurcations and a new constraint file from the multifurcated tree was constructed. A
586 second iteration of FastTree was initiated on the whole alignment including all sites to
587 produce the final tree. This tree was rooted at the outgroup.

588

589 Reconstruction of Ancestral Genome Sequences

590 Ancestral states were estimated by Fitch Traceback(15). Briefly, character sets were
591 constructed from leaf to root where each node was assigned the intersection of the
592 descendant character sets if not empty and the union otherwise. Then, moving from root
593 to leaf, nodes with more than one character in their set were assigned the consensus
594 character if present in their set or a randomly chosen representative character
595 otherwise. Substitutions between states were identified and placed in the middle of the
596 branch bridging the pair of nodes.

597 Statistical associations between mutations were computed in a manner similar to that
598 previously described(24). Briefly, sequences were leaf-weighted based on the branch
599 lengths of the ultrameterized, tree. Every mutation present across the tree at 50 mean
600 leaf-weight equivalents or more was considered. The probability of independent co-
601 occurrence between any pair was estimated in two ways. An arbitrary member of the
602 pair was selected as the ancestral mutation, and the binomial probability:

$$\sum_{k=N_{pair}}^{N_{total}} \binom{N_{total}}{k} F^k (1-F)^{N_{total}-k}$$

603

604 was computed where N_{total} is the number of substitutions to the descendant mutation
605 across the entire ancestral record, N_{pair} is the number of substitutions to the
606 descendant which succeed or appear simultaneously with a substitution to the ancestral
607 mutation, and F is the fraction of the tree (fraction of all applicable branch lengths)
608 occupied by the ancestral mutation. The ancestral/descendent designation was then
609 reversed and the “binomial score” was constructed as the negative log of the product of
610 these two terms. Additionally, for each pair, the observed and expected (product of the

611 tree fractions) tree intersections were calculated and the “Poisson score” (analogous to
612 the log-odds ratio) was calculated:

$$\begin{cases} -\ln(1 - PCDF(exp, obs)), obs > exp \\ \ln(PCDF(exp, obs)), obs < exp \end{cases}$$

613 where PCDF(exp,obs) is the cumulative probability of a Poisson distribution with mean
614 “exp”, the expected value of the data, and evaluated at “obs”, the observed value of the
615 data. Both scores are reported. Fig. 1E and Table S3 display putative positively
616 selected mutations with both scores above 5 or at least two simultaneous substitutions.
617 Fig. S10 has a relaxed score threshold for 2, an increased weight threshold of 100, and
618 is not restricted to positively selected residues.

619

620 Classical Multidimensional Scaling of the MSA

621 Pairwise Hamming distances were computed for all pairs of rows in the global MSA
622 ignoring gaps and ambiguous characters i.e. the sequences X=“ATN-A” and
623 Y=“NTAAT” would be assigned a distance of 1. The resulting distance matrix was
624 embedded in three dimensions with the MATLAB(47) routine “cmdscale”. 100 rounds of
625 stochastically initiated k-means clustering of the embedding was conducted and the
626 optimum cluster number was determined to be 4 on the basis of the silhouette score
627 distribution (Fig S1).

628

629 Validation of Mutagenic Contexts

630 Mutations were divided into four categories: synonymous vs non-synonymous
631 substitutions and high vs low frequency of independent occurrence. For example,
632 consider codon X with 3 non-synonymous substitutions gat->ggt and 1 non-synonymous
633 substitution gat->cgt. In this context, a non-synonymous nucleotide substitution a->g of
634 frequency 4 would be recorded in nucleotide (X-1)*3+2. The low vs high frequency
635 threshold was determined by the 95th percentile of the synonymous mutation frequency
636 distribution (operationally. 5). For each mutation, the trinucleotide contexts from the
637 ancestral reconstruction at the nodes where the mutation occurred were compared to
638 the background genome-wide frequencies, computed for the inferred common ancestor
639 of SARS-CoV-2.

640

641 The expected frequencies of the trinucleotides using the background distribution were
642 tabulated; the Yates correction (+/-0.5 to the original count depending on whether the
643 count is below or above the expectation) was applied to the observed frequencies; the
644 log-odds ratios of the (corrected) observed frequencies to the expectation were
645 computed; and CMDS was applied to the Euclidean distances between the log-odds

646 vectors to embed the points onto a plane (Fig. S5 A.). This analysis was then repeated,
647 this time, distinguishing only between high and low frequency substitutions but not N
648 and S (Fig. S5 B). Finally, the differences in the contexts of high frequency synonymous
649 vs non-synonymous events were considered in the same manner and the chi-square
650 statistics $((\text{observed}-\text{expected})^2/\text{expected})$ were compared with the critical chi-square
651 value ($p=0.05/64$, $df=1$, Fig. S5 C.).

652

653 Computation of dN/dS

654 For each of the 24 ORFs (splitting orf1ab into 15 segments corresponding to the 15
655 mature proteins, nsp11 and nsp12 combined), 10 reduced alignments were constructed
656 as follows. Sequences were ordered based on diversity, in the same order with which
657 they were included in the constraint trees. The first 10 sequences are conserved across
658 every alignment and the remaining 40 are unique to each alignment. The reference
659 sequence, NC_045512.2, was additionally added to each reduced alignment. PAML(48)
660 was then used to estimate tN , tS , dN/dS , N , S , and N/S for each segment and every
661 reduced alignment.

662 Given the global ancestral reconstruction from Fitch traceback, the total number of non-
663 synonymous and synonymous substitutions (nN and nS , respectively) as well as these
664 tallies normalized by the respective segment length (tN , and tS , respectively) were
665 retrieved for each segment. . A hybrid dN/dS value for each segment was estimated to
666 be $(nN/nS)/(N/S)^*$ where $(N/S)^*$ is the median value of N/S across all repeats for the
667 segment.

668

669 Metadata Assignment

670 Headers for all isolates belonging to CD-HIT clusters with a representative incorporated
671 into the alignment with fewer than 10 interior ambiguous characters were processed to
672 extract the sequencing date and location. Sequencing location abbreviations were
673 matched to full names and the latitude/longitude of a representative city for each
674 location was retrieved from simplemaps (<https://simplemaps.com/data/world-cities>)(49).

675

676 Regional Divergence Analysis

677 Two approaches, one partition dependent and one partition independent, were used.
678 First three groups of isolates were constructed belonging to partitions 1&2, partitions
679 3,4,&5, and partition 6. The Hellinger distance between regions over a sliding time
680 window was then computed between regions for this three group distribution. Next, 400
681 isolates were randomly selected from each region over a sliding window and 200 pairs
682 within each region as well as 200 pairs between each pair of regions were composed.

683 The tree distance between each pair was computed and the mean for each inter- and
684 intra-regional pair tree-distance distribution was recorded.

685

686 **Supplemental References:**

687

688 [42] W. Li, A. Godzik. Cd-hit: a fast program for clustering and comparing large sets of protein or
689 nucleotide sequences. *Bioinformatics* **22.13**, 1658-1659 (2006)

690 [43] L. Fu et al. CD-HIT: accelerated for clustering the next-generation sequencing data.
691 *Bioinformatics* **28.23**, 3150-3152 (2012)

692 [44] K. Katoh et al. MAFFT: a novel method for rapid multiple sequence alignment based on fast
693 Fourier transform. *Nucleic acids research* **30.14**, 3059-3066 (2002)

694 [45] M.N. Price, P.S. Dehal, A.P. Arkin. FastTree 2—approximately maximum-likelihood trees for
695 large alignments. *PloS one* **5.3**, e9490 (2010)

696 [46] L. Nguyen et al. IQ-TREE: a fast and effective stochastic algorithm for estimating
697 maximum-likelihood phylogenies. *Molecular biology and evolution* **32.1**, 268-274 (2015)

698 [47] MathWorks, Inc, ed. MATLAB, high-performance numeric computation and visualization
699 software: reference guide. MathWorks, (1992)

700 [48] Z. Yang. PAML 4: phylogenetic analysis by maximum likelihood. *Molecular biology and*
701 *evolution* **24.8**, 1586-1591 (2007)

702 [49] World Cities Database. *simplemaps*. <https://simplemaps.com/data/world-cities>

703

704 **Figure legends**

705

706 **Figure 1. Evolution of SARS-CoV-2.**

707 **A.** Global tree reconstruction with 6 principal partitions enumerated and color-coded.

708 The gray clades in partition 3/6 are composed of late sequences from

709 Europe/Australia&New Zealand respectively. The majority of sequences from AU/NZ

710 reside in the respective clade as do almost half of all late isolates from Europe. **B.**

711 Projections of the 3D embedding of the pairwise Hamming distance matrix between

712 SARS-CoV-2 genomes. The partitions are color-coded as in A. Wires enclose the

713 convex hulls for each of the four optimal clusters. **C.** Signatures of amino acid

714 replacements for each partition. Sites are ordered by decreasing maximum Kullback-

715 Leibler divergence of the nucleotide distribution (sites are not consecutive in the SARS-

716 CoV-2 proteins; the proteins along with nucleotide and amino acid numbers are

717 indicated underneath each column) of any site in any partition relative to the distribution

718 in that site over all partitions. **D.** Site history trees for spike 614 and nucleocapsid 203

719 positions. Nodes were included in this reduced tree based on the following criteria:

720 those immediately succeeding a substitution; representing the last common ancestor of

721 at least two substitutions; or terminal nodes representing branches of five sequences or

722 more (approximately, based on tree weight). Edges are colored according to their

723 position in the main partitions and the line type corresponds to the target mutation

724 (solid) or any other state (dashed). These sites are largely binary as are most sites in

725 the genome. The sizes of the terminal node sizes are proportional to the log of the

726 weight descendent from that node beyond which no substitutions in the site occurred.

727 Node color corresponds to target mutation (black) or any other state (gray). **E.** Network

728 of putative epistatic interactions for likely positively selected residues.

729

730 **Figure 2. Global and regional SARS-CoV-2 partition dynamics during the COVID-**

731 **19 pandemic. A.** US partition distribution over time. **B.** European partition distribution

732 over time **D.** Asian partition distribution over time. **E.** Australian/New Zealand partition

733 distribution over time. **E.** The frequency of mutation S|614G and N|194L, 119L, 203K,

734 205I, or 220V.

735

736 **Figure 3. Global and regional trends in SARS-CoV-2 evolution. A.** Global
737 distribution of sequences with sequencing locations in the US (brown), Europe (gray),
738 Asia (brown), and Australia/New Zealand (gray) identified. Color scheme is for visual
739 distinction only. Pie charts indicate the partition distributions for each region mid-March
740 through mid-April and mid-July through mid-August. **B.** The Hellinger distance between
741 the six pairs of regions over the four group distribution: partitions 1&2, partitions 3,4&5,
742 and partition 6. **C.** The ratio of the mean tree-distance for pairs of isolates between
743 regions vs. isolates within regions.

744 Supplemental Figures

745

746 **Figure S1.** 25th, median (solid line), and 75th percentiles of the silhouette score
747 distribution for 100 stochastically initiated rounds of k-means clustering for 2-10
748 clusters.

749

750 **Figure S2.** The Kullback-Leibler divergence between each clade and the whole for the
751 ten most divergent codons in the genome. The solid line indicates the maximum of any
752 clade and points represent the remaining clades.

753

754 **Figure S3. A.** Distributions of the moving average, respecting segment boundaries,
755 across a 100 codon window for synonymous (blue) and amino acid (orange)
756 substitutions. Solid lines: normal approximations of the distributions (same median and
757 interquartile distance); solid lines: approximation with the same median and theoretical
758 (Poisson) variance. **B.** Moving averages, respecting segment boundaries, across a 100
759 codon window for synonymous and nonsynonymous substitutions per site, raw (top)
760 and normalized by the median (bottom). There are several regions in the genome with
761 an apparent dramatic excess of synonymous substitutions: 5' end of orf1ab gene; most
762 of the M gene; 3'-half of the N gene, as well as amino acid substitutions: most of the
763 orf3a gene; most of the orf7a gene; most of the orf8 gene; and several regions in of the
764 N gene.

765

766 **Figure S4.** Moving average over a window of 1000 codons, not respecting segment
767 boundaries, of the total number of nucleotide exchanges $n1 \rightarrow n2$ summed over all
768 substitutions. The ratio to the median over the entire alignment is also displayed as well
769 as the normalized exchange distribution (*i.e.* $\#c \rightarrow t / (\#c \rightarrow t + \#c \rightarrow g + \#c \rightarrow a)$).

770

771 **Figure S5 A.** Two dimensional embedding of the Euclidean distances between the log-
772 odds vectors of low and high frequency, nonsynonymous and synonymous mutations in
773 the space of trinucleotide contexts relative to background expectation. The context of
774 the high-frequency events (both S and N) is dramatically different from the background
775 frequencies. There is a strong common component in the deviation of both kinds of
776 high-frequency events. The context of the low-frequency events (both S and N) also
777 differs slightly, in the same direction, from the background frequencies. There is a
778 consistent distinction between synonymous and non-synonymous events, suggesting
779 that a single mutagenic context or mechanistic bias does not account for both S and N
780 events. **B.** Log odds ratio of low and high frequency mutations, both synonymous and
781 nonsynonymous, relative to background expectation for each trinucleotide context. The

782 NCN context (i.e. all mutations C->D) harbors dramatically more mutation events than
783 the other contexts (all 16 NCN events are within the top 20 most-biased high-frequency
784 events). The log-odds ratios for low-frequency events are somewhat correlated with
785 those for high-frequency events ($r_{\text{Pearson}}=0.5$, without NCN $r_{\text{Pearson}}=0.64$),
786 suggesting the same mechanism may be responsible for the strong bias observed
787 among high frequency events and the weaker bias observed among low frequency
788 events. **C.** Log odds ratio of high frequency nonsynonymous mutations relative to the
789 background expectation from the sum of both high synonymous and high
790 nonsynonymous mutations vs. the sum + 1. There are 12 contexts where synonymous
791 and non-synonymous events differ significantly. All contexts with an excess of
792 synonymous events are NCN, suggesting that high-frequency synonymous events
793 could be driven by mechanistic bias; on the contrary, none of the contexts with an
794 excess of non-synonymous mutations are NCN (in fact all are NGN:
795 aga,agt,agg,ggt,agc,tgt,gct vs. aca,tct,tca,tcg,ccg), suggesting that these non-
796 synonymous events could be driven by other mechanisms. There is no correlation
797 between the frequency of event context and the log-odds ratio for non-synonymous
798 events, further suggesting that the log-odds ratio is not biased by hot-spot mutation
799 context.

800

801 **Figure S6.** Correspondence between the “tree length for dN”, “tree length for dS”, and
802 dN/dS between PAML and the results of the ancestral reconstruction utilizing Fitch
803 traceback across 24 ORFs. Two outliers in the PAML tS distribution are identified in
804 each plot.

805

806 **Figure S7. A.** The number of nonsynonymous events vs the number of synonymous
807 events per codon. **B.** The moving average of 100 codons, respecting segment
808 boundaries. **C.** The moving average after removing events with 5 or more independent
809 occurrences. Rho refers to Spearman. Dashed lines are $2/1.3^*x$ reflecting the genome-
810 wide ratio of nonsynonymous to synonymous substitutions, solid lines are linear best fit.

811

812 **Figure S8.** The fraction of sites with at least one substitution vs moving averages,
813 respecting segment boundaries, over windows of 100 codons for synonymous and
814 nonsynonymous substitutions.

815

816 **Figure S9.** Structural analysis for sites epistatically linked to spike D614 within the spike
817 protein. D614 is at the interface between Spike chains. Most regions in the vicinity are
818 not structurally solved potentially indicating that depending on the status of the RBD of
819 the other chains, the regions in close proximity to D614 could become highly flexible.
820 Residue 21 is not structurally solved; however, model inference suggests it is spatially

821 distant from residue 614. H49 makes a stack cation pi interaction with R44 within the
822 same chain. H49 is spatially distant from D614, however, the domain it belongs to
823 (circled in red) is linked by a linker (dashed red line) that leads to the domain containing
824 D614 (circled in purple). This potentially functions as a holding point to position the
825 purple domain. Note that 614 is very close to the cleavage site, likely requiring accurate
826 positioning of this domain.

827

828 **Figure S10.** Epistatic network for the tree including mutations with binomial/poisson
829 scores above 2 or at least two simultaneous substitutions and weight of at least
830 approximately 100 leaves not restricted to likely positively selected residues.

831

832 **Figures S11.** Correlation between sequencing date and tree distance to the root for all
833 isolates with metadata as well as those which appear explicitly in the tree.

834

835 **Figures S12-13.** Global distribution of sequences. Color represents the number of
836 sequences from that location and size represents the fraction of sequences from the
837 clade displayed. Clade indices are in the top left corner of each map.

838

839 **Figure S14.** Major partition distributions for each region at two fixed timepoints, mid-
840 March to mid-April and mid-July to mid-August, as well as the difference.

841

842 **Figure S15.** The Kullback-Leibler divergence between the gray clade (solid line) and
843 remaining sequences (dashed line) vs. the whole of partition 6. The ten most divergent
844 codons in the genome are shown in the sequence logo.

845

846 **Figure S16.** The Kullback-Leibler divergence between the gray clade (solid line) and
847 remaining sequences (dashed line) vs. the whole of partition 3. The ten most divergent
848 codons in the genome are shown in the sequence logo.

849

850 **Figure S17.** The sequence logo for the 15 most divergent (Kullback-Leibler) codons in
851 the genome considering US sequences sourced after July 17 vs. all sequences sourced
852 after July 17.

853

854 **Figure S18.** The sequence logo for the 15 most divergent (Kullback-Leibler) codons in
855 the nucleocapsid protein considering US sequences sourced after July 17th 2020 vs. all
856 sequences sourced after July 17th 2020.

857

858 **Figure S19.** The frequencies of mutations S|614G, N|194L, N119L, N203K, N205I, and
859 N220V over time.

860

861 **Figure S20.** The mean tree distance between pairs of isolates **A.** from different regions,
862 **B.** within the same region, and **C.** the ratio over time.

863

864 **Supplemental Tables**

865

866 **Table S1.** The top ten mutations most commonly observed and the top ten with the
867 greatest number of parallel substitutions.

868

869 **Table S2.** List of sites likely to be evolving under positive selection.

870

871 **Table S3.** All epistatic interactions among states meeting the criteria outlined in the
872 main text for likely positive selection with binomial/Poisson scores greater than 5 or at
873 least 2 simultaneous substitutions. Each mutation must have a minimum weight of
874 approximately 50 leaves and each pair, 30 leaves. Each pair is arbitrarily ordered and
875 the numbers of simultaneous, descendant, and independent substitutions are tabulated.

876

877 **Table S4.** Tabulated three codon neighborhoods for all sites containing at least one
878 stop codon. Sites are ordered in decreasing number of sequences containing the stop.
879 Stops are listed separately before all other neighborhoods.

880

881 **Table S5.** The number of isolates (out of approximately 73k) observed to bear at least
882 one substitution relative to the reference sequence, NC_045512.2, within the regions
883 specified. These regions are commonly used within PCR assays for diagnostic testing.

884

885





