1

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

4

Nash D. Rochman[1,*], Yuri I. Wolf[1], Guilhem Faure[2], Pascal Mutz[1], Feng Zhang[2,3,4,5,6,*] and Eugene V. Koonin[1,*]

[1]National Center for Biotechnology Information, National Library of Medicine, Bethesda, MD 20894
[2]Broad Institute of MIT and Harvard, Cambridge, MA 02142; [3]Howard Hughes Medical Institute, Massachusetts Institute of Technology, Cambridge, MA 02139; [4]McGovern Institute for Brain Research, Massachusetts Institute of Technology, Cambridge, MA 02139; [5]Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, MA 02139; and [6]Department of Biological Engineering, Massachusetts Institute of Technology, Cambridge, MA 02139

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

15

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

17

18

## Abstract

Understanding the trends in SARS-CoV-2 evolution is paramount to control the COVID-19 pandemic. We analyzed more than 300,000 high quality genome sequences of SARS-CoV-2 variants available as of January 2021. The results show that the ongoing evolution of SARS-CoV-2 during the pandemic is characterized primarily by purifying selection, but a small set of sites appear to evolve under positive selection. The receptor-binding domain of the spike protein and the nuclear localization signal (NLS) associated region of the nucleocapsid protein are enriched with positively selected amino acid replacements. These replacements form a strongly connected network of apparent epistatic interactions and are signatures of major partitions in the SARS-CoV-2 phylogeny. Virus diversity within each geographic region has been steadily growing for the entirety of the pandemic, but analysis of the phylogenetic distances between pairs of regions reveals four distinct periods based on global partitioning of the tree and the emergence of key mutations. The initial period of rapid diversification into region-specific phylogenies that ended in February 2020 was followed by a major extinction event and global homogenization concomitant with the spread of D614G in the spike protein, ending in March 2020. The NLS associated variants across multiple partitions rose to global prominence in March-July, during a period of stasis in terms of inter-regional diversity. Finally, beginning July 2020, multiple mutations, some of which have since been demonstrated to enable antibody evasion, began to emerge associated with ongoing regional diversification, which might be indicative of speciation.

## Significance

Understanding the ongoing evolution of SARS-CoV-2 is essential to control and ultimately end the pandemic. We analyzed more than 300,000 SARS-CoV-2 genomes available as of January 2021 and demonstrate adaptive evolution of the virus that affects, primarily, multiple sites in the spike and nucleocapsid protein. Selection appears to act on combinations of mutations in these and other SARS-CoV-2 genes. Evolution of the virus is accompanied by ongoing adaptive diversification within and between

49    geographic regions.  This diversification could substantially prolong the pandemic and

50    the vaccination campaign, in which variant-specific vaccines are likely to be required.

51

52

## Introduction

High mutation rates of RNA viruses enable adaptation to hosts at a staggering pace (1-4). Nevertheless, robust sequence conservation indicates that purifying selection is the principal force shaping the evolution of virus populations, with positive selection affecting only relatively small subsets of sites directly involved in virus-host coevolution (5-8). The fate of a novel zoonotic virus is in part determined by the race between public health intervention and virus diversification. Even intermittent periods of positive selection can result in lasting immune evasion, leading to oscillations in the size of the susceptible population, and ultimately, a regular pattern of repeating epidemics, as has been amply demonstrated for Influenza(9-11).

During the current coronavirus pandemic (COVID-19), understanding the degree and dynamics of the diversification of severe acute respiratory syndrome coronavirus 2 (SARS-Cov-2) and identification of sites subject to positive selection are essential for establishing practicable, proportionate public health responses, from guidelines on isolation and quarantine to vaccination(12).  To investigate the evolution of SARS-CoV-2, we collected all available SARS-Cov-2 genomes as of January 8, 2021, and constructed a global phylogenetic tree using a "divide and conquer" approach. Patterns of repeated mutations fixed along the tree were analyzed in order to identify the sites subject to positive selection. These sites form a network of potential epistatic interactions. Analysis of the putative adaptive mutations provides for the identification of signatures of evolutionary partitions of SARS-CoV-2. The dynamics of these partitions over the course of the pandemic reveals alternating periods of globalization and regional diversification.

## Results and Discussion

### Global multiple sequence alignment of the SARS-CoV-2 genomes

4

83  To investigate the  evolution of SARS-CoV-2, we aggregated all available SARS-Cov-2

84  genomes as of January 8, 2021, from the three principal repositories: Genbank(13),

85  Gisaid(14), and CNCB(15). From the total of 321,096 submissions in these databases,

86  175,857 unique SARS-Cov-2 genome sequences were identified, and 98,185 high

87  quality sequences were incorporated into a global multisequence alignment (MSA)

88  consisting of the concatenated open reading frames with stop codons trimmed. The vast

89  majority of the sequences excluded from the MSA were removed due to a

90  preponderance of ambiguous characters (see Methods). The sequences in the final

91  MSA correspond to 175,776 isolates with associated date and location metadata.

92

93  **Tree Construction**

94

95  Several methods for coronavirus phylogenetic tree inference have been tested(16, 17).

96  The construction of a single high-quality tree from nearly 200,000 30 kilobase (kb)

97  sequences using any of the existing advanced methods is computationally prohibitive.

98  Iterative construction of the complete phylogeny would seem an obvious solution such

99  that a global topology would be obtained based on a subset of sequences available at

100 an earlier date, and later sequences would be incorporated into the existing tree.

101 However, this approach induces artifacts through the inheritance of deep topologies that

102 differ substantially from any maximum-likelihood solution corresponding to the complete

103 alignment.

104

105 Therefore, building on the available techniques, we utilized a "divide and conquer"

106 approach which is not subject to these artifacts and furthermore can be employed for

107 datasets that cannot be structured by sequencing date, including metagenomic

108 analyses. This approach leverages two ideas. First, for any alignment, a diverse

109 representative subset of sequences can be used to establish a deep topology, the tree

110 "skeleton", that corresponds to a maximum-likelihood solution over the entire alignment.

111 Second, deep branches in an unrooted tree are primarily determined by common

112 substitutions relative to consensus. In other words, rare substitutions are unlikely to

5

113 affect deep splits or branch lengths. We adopted the following steps to resolve the

114 global phylogeny for SARS-CoV-2 (see Methods for details).

115

116 The first step is to construct sets of diverse representative sequences such that the

117 topologies inferred from each subset share the same tree "skeleton" corresponding to

118 that of the global topology. Sequence diversity is measured by the hamming distance

119 between pairs of sequences; however, maximizing hamming distances among a set of

120 representative sequences does not guarantee maximization of the tree distances in the

121 resultant global topology and so does not guarantee the maximum-likelihood topology of

122 this subset would share the global tree skeleton. Therefore, a reduced alignment

123 containing only the top 5% of sites (ignoring nearly-invariant sites, see below) with the

124 most common substitutions relative to consensus was constructed. Sequences

125 redundant over this narrow alignment were removed. Sets of diverse, based on the

126 hamming distances over this reduced alignment, representative sequences were then

127 achieved and all subtrees generated from each diverse subset were aggregated to

128 constrain a single, composite tree.

129

130 This composite tree reflects the correct tree skeleton and could be used to constrain the

131 global topology; however, due to numerous sequencing errors in this dataset, another

132 intermediate step was taken. A second reduced alignment was constructed in which

133 nearly-invariant sites, which may represent sequencing errors and should not be used

134 to infer tree topology, were omitted. As before, sequences redundant over this reduced

135 alignment were removed and a tree was then constructed from this alignment,

136 constrained to maintain the topology of the composite tree wherever possible. This tree

137 reflects the correct topology of the global tree but has incorrect branch lengths. Finally,

138 the global tree with the correct branch lengths (Fig. 1A) was constructed over the whole

139 alignment, constrained to maintain the topology of previous tree wherever possible. A

140 complete reconstruction of ancestral sequences was then performed by leveraging Fitch

141 traceback(18), enabling comprehensive identification of nucleotide and amino acid

142 replacements across the tree.

143

6

144    We identified 8 principal partitions within this tree, in a general agreement with other

145    work(19-21), along with three divergent clades (Fig. 1A) that, as discussed below, are

146    important for the interpretation of the metadata. Given the short evolutionary distances

147    between SARS-CoV-2 isolates, despite the efforts described above, the topology of the

148    global tree is a cause of legitimate concern(17, 22-24). For the analyses presented

149    below, we rely on a single, explicit tree topology which is probably one of many equally

150    likely estimates(17). Therefore, we sought to validate the robustness of the major

151    partitions of the virus genomes using a phylogeny-free approach. To this end, pairwise

152    Hamming distances were computed for all sequences in the MSA and the resulting

153    distance matrix was embedded in a 3-dimensional subspace using classical

154    multidimensional scaling. In this embedding, the 8 partitions are separated, and the

155    optimal clustering, determined by *k*-means, returned 5 categories (see Methods, Fig.

156    S1), of which 4 correspond to partitions 5 and 8, and the divergent clades v1 and v2.

157    These findings indicate that an alternative tree with a comparable likelihood but a

158    dramatically different coarse-grain topology, most likely, cannot be constructed from this

159    MSA.

160

161    **Mutational Signatures&Biases and Estimation of Selection**

162

163    Each of the 8 partitions and 3 variant clades can be characterized by a specific amino

164    acid replacement signature (Fig. 1B), generally, corresponding to the most prominent

165    amino acid replacements across the tree (Table S1), some of which are shared by two

166    or more partitions and appear independently many times, consistent with other

167    reports(25). The receptor binding domain (RBD) of the spike protein and a region of the

168    nucleocapsid protein associated with nuclear localization signals (NLS)(26) are enriched

169    with these signature replacements, but they are also found in the nonstructural proteins

170    1ab, 3a, and 8. The identification of these prevailing non-synonymous substitutions and

171    an additional set of frequent synonymous substitutions suggested that certain sites in

172    the SARS-CoV-2 genome might be evolving under positive selection. However,

173    uncovering the selective pressures affecting virus evolution was complicated by non-

174    negligible mutational biases. The distributions of the numbers of both synonymous and

175   non-synonymous substitutions across the genome were found to be substantially over-

176   dispersed compared to both the Poisson and normal expectations (Fig. S2).

177   Examination of the relative frequencies of all 12 possible nucleotide substitutions

178   indicated a significant genome-wide excess of C to U mutations, approximately

179   threefold higher than any other nucleotide substitution, with the exception of G to U, as

180   well as some region-specific trends. Specifically, G to U mutations increase steadily in

181   frequency throughout the second half of the genome and the distribution of nucleotide

182   substitutions over the polyprotein is dramatically different from other ORFs (Fig. S3).

183

184   Motivated by the observation of the mutational biases, we compared the trinucleotide

185   contexts of synonymous and non-synonymous substitutions as well as the contexts of

186   low and high frequency substitutions. The contexts of high-frequency events, both

187   synonymous and non-synonymous, were found to be dramatically different from the

188   background frequencies. The NCN context (that is, all C->D mutations) harbors

189   substantially more events than other contexts (all 16 NCN triplets are within the top 20

190   most high-frequency-biased, see Methods and Fig. S4) and is enriched in mutations

191   uniformly across the genome, primarily, among high-frequency sites. This pattern

192   suggests a mechanistic bias of the errors made by the coronavirus RNA-dependent

193   RNA polymerase (RdRP). Evidently, such a bias that increases the likelihood of

194   observing multiple, independent substitutions in the NCN context complicates the

195   detection of selection pressures. However, only 2 of the 9 contexts with an excess of

196   non-synonymous events are NCN (gct,tct, Fig. S4), suggesting that at least some of

197   these repeated, non-synonymous mutations are driven by other mechanisms. Thus, we

198   excluded all synonymous substitutions and non-synonymous substitutions with the NCN

199   context from further consideration in the determination of candidate sites evolving under

200   positive selection.

201

202   Beyond this specific context, the presence of any hypervariable sites complicates the

203   computation of the *dN/dS* ratio, the gauge of protein-level selection(27), which requires

204   enumerating the number of synonymous and non-synonymous substitutions within each

205   gene. Hypervariable sites bias this analysis, and therefore, we used two methods to

8

206   ensure reliable estimation of *dN/dS*. For each protein-coding gene of SARS-CoV-2

207   (splitting the long orf1ab into 15 constituent non-structural proteins), we obtained both a

208   maximum likelihood estimate of *dN/dS* across 10 sub-alignments and an approximation

209   computed from the global ancestral reconstruction (see Methods). This approach was

210   required due to the size of the alignment, which makes a global maximum likelihood

211   estimation computationally prohibitive. Despite considerable variability among the

212   genes, we obtained estimates of substantial purifying selection (0.1<*dN/dS*<0.5) across

213   most of the genome(Fig. S5), with a reasonable agreement between the two methods.

214   This estimate is compatible with previous demonstrations of purifying selection affecting

215   about 50% of the sites surveyed or more(5),  among diverse RNA viruses(6).

216

217   **Evidence of Positive Selection**

218   As shown in the previous section, evolution of SARS-CoV-2 is likely primarily driven by

219   substantial purifying selection. However, more than 100 non-synonymous substitutions

220   appeared to have emerged multiple times, independently, covering a substantial portion

221   of the tree equivalent to approximately 200 or more terminal branches or "leaves", and

222   were not subject to an overt mechanistic bias. Due to the existence of many equally

223   likely trees, in principle, in one or more of such trees, any of these mutations could

224   resolve to a single event. However, such a resolution would be at the cost of inducing

225   multiple parallel substitutions for other mutations, and thus, we conclude that more than

226   100 codons in the genome that are not subject to an overt mechanistic bias underwent

227   multiple parallel mutations in the course of SARS-CoV-2 evolution during the COVID-19

228   pandemic.

229

230   One immediate explanation of this observation is that these sites evolve under positive

231   selection. The possible alternatives could be that these sites are mutational hotspots or

232   that the appearance of multiple parallel mutations was caused by numerous

233   recombination events (either real or artifacts caused by incorrect genome assembly

234   from mixed infections) in the respective genomic regions. Contrary to what one would

235   expect under the hotspot scenario, we found that codons with many synonymous

236   substitutions tend to harbor few non-synonymous substitutions, and vice versa (Fig. S6

237  A). When a moving average with increasing window size was computed, only a weak

238  positive correlation was observed between the numbers of synonymous and non-

239  synonymous substitutions (Figs. S6 B&C, S7). Most sites in the virus genome are highly

240  conserved, the sites with most substitutions tend to reside in conserved neighborhoods,

241  and the local fraction of sites that harbor at least one mutation strongly correlates with

242  the moving average (Fig. S8). Together, these observations indicate that SARS-CoV-2

243  genomes are subject to diverse site-specific and regional selection pressures but we did

244  not detect regions of substantially elevated mutation or recombination in general

245  agreement with other studies(28) despite the role recombination might have played in

246  zoonosis(29-34).

247

248  **Positively selected sites in SARS-CoV-2 proteins**

249

250  Given the widespread purifying selection affecting evolving SARS-CoV-2 genomes,

251  substantially relaxed selection at any site is expected to permit multiple, parallel non-

252  synonymous mutations to the same degree that any site harbors multiple, parallel

253  synonymous mutations. Thus, seeking to identify sites subject to positive selection, we

254  focused only on those non-synonymous substitutions that independently occurred more

255  frequently than 90% of all synonymous substitutions excluding the mutagenic NCN

256  context (see Methods). Most if not all sites in the SARS-CoV-2 genome that we found to

257  harbor such frequent, parallel non-synonymous substitutions outside of the NCN context

258  can be inferred to evolve under positive selection (Table S2, List 1). The positively

259  selected residues form a co-occurrence network that likely reflects epistatic interactions

260  (Fig. 1D and Table S3, see Methods), in which the central hubs are D614G in the spike

261  (S) protein and two adjacent substitutions in the nucleocapsid (N) protein, R203K and

262  G204R, the three most common positively selected mutations (Fig. 1C) (35). Fig. 1D,).

263

264  **Positively selected amino acid replacements in the receptor-binding domain of**

265  **the spike protein**

266

267  Spike D614G appears to boost the infectivity of the virus, possibly, by increasing the

268  binding affinity between the spike protein and the cell surface receptor of SARS-CoV-2,

269  ACE2(36). Conclusively demonstrating selection for a single site has proven

270  challenging(37), even within this robust dataset. Although the emergence of this

271  mutation corresponds to the extinction of partitions lacking 614G (see below), the

272  possibility remains that this mutation is a passenger to some other mutagenic or

273  epidemiological event. The 614 site of the S protein is evolutionarily labile, so that  the

274  ancestral reconstruction includes multiple gains of 614D after a previous loss. As a

275  result,  the reverse replacement G614D appears often enough to pass our statistical

276  criteria for positive selection. Although severely biasing against recent events, one can

277  additionally require that the mean tree fraction descendant from each candidate

278  positively selected amino acid replacement be sufficiently large, removing from

279  consideration events which are frequent but shallow (see Methods). The addition of this

280  criterion results in a "shortlist" of 22 residues subject to the strongest selection (Table

281  S2, List 2) that do not include 614D.

282

283  Additionally, apart from the selective advantage of a single replacement, it should be

284  emphasized that D614G (but not G614D) is a central hub of the epistatic network (Fig.

285  1D). Conceivably, epistatic interactions with this residue can result in ensembles of

286  mutations which substantially increase fitness. The ubiquitous epistasis throughout

287  molecular evolution(38-41) suggests the possibility that many if not most mutations,

288  which confer a substantial selective advantage, do so only within a broader epistatic

289  context, not in isolation. By increasing the receptor affinity, D614G apparently opens up

290  new adaptive routes for later steps in the viral lifecycle. The specific mechanisms of

291  such hypothetical enhancement of virus reproduction remain to be investigated

292  experimentally.

293

294  In addition to 614G, 31 spike mutations, most within the RBD, are signature mutations

295  for divergent clades v1-3; emergent variants vAfrica or vOceania (see below); or

296  established variants B.1.1.7, B.1.1.7_E484K, B.1.258_delta, B.1.351, B.1.429, P.1, or

297  P.2(42-47) (Table S4, List 1). Three of these signature mutations pass the strict criteria

11

298    for positive selection: S|N501Y, S|S477N, and S|V1176F, and S|N501Y makes the

299    shortlist of the 22 strongest candidates. H69del/V70del are signature mutations for

300    variant B.1.258_delta and have been previously observed to have rapidly emerged in an

301    outbreak among minks(48, 49). A two amino acid deletion (in our alignment this deletion

302    resolves to sites 68/69 due to many ambiguous characters in this neighborhood)

303    appears multiple times independently throughout the tree and is present in

304    approximately one third of the European sequences from January, 2021(Fig. S9,

305    deletions are not shown in Fig. 1B, see below).

306

307    Two sites within the RBD, N331 and N343, have been shown to be important for the

308    maintenance of infectivity(50). As could be expected, these amino acid residues are

309    invariant. Four more substitutions in the RBD, among others, N234Q, L452R, A475V,

310    and V483A, have been demonstrated to confer antibody resistance(50). N234Q, A475V,

311    and V483A were never or rarely found in our alignment but L452R is a signature of

312    variant B.1.429. Although not meeting our criteria for positive selection, it appeared

313    multiple times across the tree, including within partition 1. Of greatest concern is

314    perhaps N501Y. This amino acid replacement is a signature of variants B.1.1.7,

315    B.1.1.7_E484K, B.1.351, P.1; divergent clade v2; and emergent variant vAfrica. N501Y

316    is among the 22 strongest candidates for positive selection and has been demonstrated

317    to escape neutralizing antibodies(51). N501T in the same site is of additional

318    concern(52) and has also been observed in mink populations(53). Additionally,

319    S|N439K, a signature mutation for variant B.1.258_DELTA that has been demonstrated

320    to enable immune escape(54), is observed in a large portion of the tree.

321

322    The emergence of multiple mutations associated with immune evasion during a period

323    of the pandemic when the majority of the global population had remained naïve is

324    striking. Such adaptations are generally expected to emerge among host populations

325    where many individuals have acquired immunity either through prior exposure or

326    vaccination(55-58). Furthermore, this pattern of, most likely independent, emergence of

327    persisting variants among both human and mink populations suggests the possibility

328    that these mutations represent non-specific adaptations acquired shortly after zoonosis.

12

329   The factors underpinning the evolution of viral life history traits after zoonosis, especially

330   virulence, remain poorly understood(59) but apparently result from selective pressures

331   imposed by both epidemiological parameters (host behavior)(60, 61), which may be

332   conserved across a variety of novel hosts, and specific properties of the host receptor.

333   Whereas emergent mutations in the RBD of SARS-CoV-2 are, for obvious reasons,

334   surveyed with great intensity, we have to emphasize the enrichment of positively

335   selected residues in the N protein, which might relate to more deeply taxonomically

336   conserved routes of host adaptation for beta-coronaviruses.

337

338   **Amino acid replacements associated with the nuclear localization signals in the**

339   **nucleocapsid protein**

340

341   Evolution of beta-coronaviruses with high case fatality rates including SARS-CoV-2 was

342   accompanied by accumulation of positive charges in the N protein that might enhance

343   its transport to the nucleus(62).  Thirteen amino acid replacements in the N protein are

344   signatures among the variants or major partitions discussed here, 7 of which: 203K,

345   204R, 205I, 206F, 220V, 234I, and 235F, are in the vicinity of the known NLS motifs or

346   other regions responsible for nuclear shuttling(26). Two additional substitutions, 194L

347   and 199L, rose to prominence in multiple regions during the summer of 2020. Two of

348   these NLS-adjacent amino acid replacements, R(agg)203K(aaa) and G(gga)204R(cga),

349   almost always appear together. This pair of substitutions includes the second and third

350   most common positively selected sites after S614,  and although another adjacent site,

351   S(agt)202N(aat) is not a signature mutation, it is the 8[th] most common positively

352   selected residue. Among the 22 nonsynonymous substitutions that are apparently

353   subject to the strongest selection (Table S2, List 2), 6 are in the N protein (202N, 203K,

354   204R, 234I, 292T, and 376T).

355

356   The replacements R(agg)203K(aaa) and G(gga)204R(cga) occur via three adjacent

357   nucleotide substitutions. R(agg)203K(aaa) resolves to two independent mutations in the

358   ancestral reconstruction: first, R(agg)203K(aag), then K(aag)203K(aaa). Furthermore,

359   the rapid rise of 220V (excluded from consideration as a candidate for positive selection

13

360    in our analysis due to its NCN context) in a European cohort during the summer of 2020

361    might be related to a transmission advantage of the variant carrying this

362    substitution(63).These substitutions, in particular G(gga)204R(cga), which increases the

363    positive charge, might contribute to the nuclear localization of the N protein as well. This

364    highly unusual cluster of multiple signature and positively selected mutations across 5

365    adjacent residues in the N protein is a strong candidate for experimental study that

366    could illuminate the evolution and perhaps the mechanisms of SARS-CoV-2

367    pathogenicity.

368

369    In addition to the many mutations of interest in the N and S proteins, Orf3a|Q57H is a

370    signature mutation for partitions 6, 7, and v1. Q57H is the 4[th] most common positively

371    selected mutation. Although not considered a candidate for positive selection in our

372    analysis due to its NCN context, ORF8 S84L is a hub in the larger epistatic network

373    including all strongly associated residues (Fig. S10).

374

375    We also identified numerous nonsense mutations. Of particular interest seems to be

376    ORF8|Q27*, which is a signature for variants B.1.1.7 and B.1.1.7_E484K and could be

377    epistatically linked to positively selected residues including N|R203K and S|D614G.

378    ORF8 has been implicated in the modulation of host immunity by SARS-CoV-2, so

379    these truncations might play a role in immune evasion(64, 65).

380

381    **Potential role of epistasis in the evolution of SARS-CoV-2**

382

383    Epistasis in RNA virus evolution, as demonstrated for influenza, can constrain the

384    evolutionary landscape and promote compensatory variation in coupled sites, providing

385    an adaptive advantage which would otherwise impose a prohibitive fitness cost(66-68).

386    Because even sites subject to purifying selection can play an adaptive role through

387    interactions with other residues in the epistatic network(69), the networks presented

388    here (Figs. 1D, S10) likely underrepresent the extent of epistatic interactions occurring

389    during SARS-CoV-2 evolution. The early evolutionary events that shaped the epistatic

390    network likely laid the foundation for the diversification of the virus relevant to virulence,

14

391   immune evasion, and transmission. As discussed above, these early mutations

392   (including S|G614D) might provide only a modest selective advantage in isolation but

393   exert a much greater effect through multiple epistatic interactions.

394

395   The epistatic network will continue to evolve through the entirety of the pandemic, and

396   indeed, all emerging variants at the time of this writing are defined not by a single

397   mutation but by an ensemble of signature mutations. Moreover, in addition to the

398   apparent widespread intra-protein epistasis, there seem to exist multiple epistatic

399   interactions between the N and S proteins. In particular, S|N501Y and N|S235F are both

400   signature mutations for variants B.1.1.7 and B.1.1.7_E484K (Table S4, List 2) and this

401   pair is in the top 25% of co-occurring pairs in our network ranked by lowest probability of

402   random co-occurrence.

403

404   As with early founder mutations, when a new variant emerges with multiple signature

405   mutations, it is unclear which, if any, confer a fitness advantage. Although it is natural to

406   focus on substitutions within the RBD, we emphasize that all emergent variants contain

407   substitutions in in the vicinity of known NLS motifs. In fact, the most statistically

408   significant signature mutation (based on the Kullback-Leibler divergence) for vAfrica

409   (consistent with variant B.1.351, see below) is N|T205I. As we suggest for S|D614G,

410   these variant signature mutations are likely to exert a greater influence through multiple

411   epistatic interactions than in isolation and each signature mutation can be a member of

412   multiple epistatic ensembles beyond the group of signature mutations within which it

413   was originally identified. Indeed signature mutations are shared among defined variants

414   and we find evidence for an additional 18 putative epistatic interactions between variant

415   signature mutations and other events throughout the tree which are not identified as

416   signature mutations for any defined variant (Table S4, List 3). The growing ensemble of

417   signature mutations that appear to be subject to positive selection and the existence of

418   a robust network of putative epistatic interactions including these signatures, suggest

419   that ongoing virus diversification is driven by host adaptation rather than occurring

420   simply by neutral drift.

421

15

## Epidemiological Trends and Ongoing Diversification of SARS-CoV-2

Analysis of within-patient genetic diversity of SARS-CoV-2 has shown that the most common mutations are highly diverse within individuals(70-72). Such diversity could either result from multiple infections, or otherwise, could point to an even greater role of positive selection affecting a larger number of sites than inferred from our tree. Similarly to the case of Influenza, positive selection on these sites could drive virus diversification and might support a regular pattern of repeat epidemics, with grave implications for public health. An analysis of the relationships between the sequencing date and location of each isolate and its position within the tree can determine whether diversification is already apparent within the evolutionary history of SARS-CoV-2.

We first demonstrated a strong correlation between the sequencing date of SARS-CoV-2 genomes and the distance to the tree root (Fig. S11), indicating a sufficiently low level of noise in the data for subsequent analyses. Examination of the global distribution of each of the major SARS-CoV-2 partitions (Figs. S12-14) indicates dramatic regional differences and distinct temporal dynamics (Fig. 2). A measure of virus diversity is necessary to map to these trends. We considered two modes of diversity. Intra-regional diversity reflects the mutational repertoire of the virus circulating in any individual region within any window of time. To measure intra-regional diversity, we sampled pairs of isolates from each region and timepoint and computed the mean tree-distance for a representative ensemble of these pairs. We found that intra-regional diversity has been steadily increasing throughout the entirety of the pandemic, with the exception of Oceania from June-August, 2020 (Figs. 3A/B) which corresponds to the period following a bottleneck in the total number of infections (Fig. S15) within that region. This unabated intra-regional diversification is a further evidence of a large repertoire of host-adaptive mutations of SARS-CoV-2 evolving within the human population.

The inter-regional diversity measures the degree to which the virus can be categorized into region-specific subtypes. A demonstration of substantial inter-regional diversity would perhaps constitute the most compelling and concerning evidence of the potential

16

453    for repeat epidemics. We developed two measures of inter-regional diversity. The first

454    one is analogous to the intra-regional diversity measure. We sampled pairs of isolates

455    within each region and between each pair of regions within the same time window,

456    computed the mean tree-distance for both representative ensembles of these pairs

457    (intra- and inter-regional pairs), and calculated the ratio of inter-regional and intra-

458    regional values (Fig. S16). The second one is a partition-level measure. For every pair

459    of regions over each time window, we computed the Hellinger distance of the 11-group

460    frequency distribution between all pairs of regions over each time window. (See

461    Methods for details).

462

463    Both measures of inter-regional diversity support the division of the pandemic, through

464    the beginning of 2021, into four periods (Fig1. 3C). The first period that ended in

465    February 2020 represents rapid diversification into region-specific phylogenies. This

466    period was followed by a major extinction event and global homogenization ending in

467    March 2020. The following five months, March-July, represented a period of stasis, in

468    terms of inter-regional diversity. Finally, July 2020 was the start of the ongoing period of

469    inter-regional diversification.

470

471    The extinction of the earliest partitions, 1 and 2, corresponds to the advent of S|D614G,

472    which became fixed in all other partitions and was globally ubiquitous by June 2020

473    (Fig. 3D). Partition 8, the only partition where N|203K and 204R were fixed, became

474    dominant in every region outside of North America in the period that followed (Fig. S17).

475    However, this did not result in a global selective sweep that would involve the extinction

476    of partitions 1-7. Instead, multiple NLS-associated mutations rose to prominence across

477    different partitions, becoming globally dominant by September (Figs. 3D, S18). To

478    resolve this trend, at least two principal variants, N|203K/204R in partition 8 and N|220V

479    in partition 5, have to be considered, and we identified 6 key amino acid replacements

480    of interest for this period (N|203K/204R, N|220V, N|199L, N|194L, N205I, N206F).

481

482    In the next phase of the pandemic, partition 8 dramatically fell from dominance in two

483    regions, Africa and Oceania, replaced by partitions 6 and 7. Although we did not find a

17

484 distinct mutational signature associated with the rise of partition 7 in Oceania (Fig. S19),

485 signatures associated with the rise of partition 6 were identified in both regions (Figs.

486 S20-21). Neither of these two groups of sequences (late sequences from partition 6,

487 Oceania and Africa, respectively) form topologically distinct clades; however, due to the

488 conserved mutational signatures, we considered both groups to represent distinct

489 emerging variants, vOceania and vAfrica. The signature for vAfrica is consistent with

490 variant B.1.351. Additionally, two divergent clades within partition 8 and one clade within

491 partition 3 emerged.

492

493 The most prominent is clade v2 with a signature consistent with variant B.1.1.7.

494 Altogether, resolving this trend of emerging substitutions in the RBD (Figs. 3D, S22-24)

495 requires the consideration of at least 3 variants and includes 59 signature mutations.

496 Clade v1 appeared first in Europe in April, 2020, v2, also in Europe, in September,

497 2020, and V3 in Asia and North America, in April, 2020 (Fig. S25). Also notably,

498 although S|477N initially appears in February/March, 2020 in Europe, Oceania, and

499 North America, it dramatically rises to prominence in Oceania in April, about 3 months

500 before this mutation becomes prominent elsewhere. S|477N is a signature mutation for

501 v1 stemming from partition 3; however, the sequences from Oceania bearing this

502 mutation from Summer, 2020 are in partition 8. The dramatic diversity of signature

503 mutations among these variants decreases the likelihood of future selective sweeps (in

504 the absence of bottlenecks in the total number of infected hosts) and increases the

505 likelihood of repeat epidemics.

506

507 **The impact of SARS-CoV-2 Diversification on Testing and Vaccination**

508

509 The ongoing diversification of SARS-CoV2 poses problems for both testing and

510 vaccination. Substitutions in the E protein have already been demonstrated to interfere

511 with a common PCR assay(73). Generally, ORF1ab is more conserved than the S

512 protein, which itself is more conserved than the remaining ORFs (Figs. S2-3). Using our

513 SARS-CoV-2 MSA, we surveyed 10 regions from ORF1ab(5), N(4), and E(1) genes that

514 are commonly used in PCR assays(74) for substitutions relative to the reference

18

515   sequence. Among the more than 175k genome sequences, there were thousands of

516   nucleotide substitutions in each of these regions, but those in ORF1ab were markedly

517   less variable than those in N (Supplementary table 5), with one region in N

518   demonstrating variability in nearly one third of all isolates. It can be expected that most

519   targets within the polyprotein will remain subject to the fewest polymorphism-induced

520   false negatives even as the virus continues to diversify.

521

522   Of the 9 primary vaccines/candidates (75), three are inactivated whole-virus (Sinovac,

523   Wuhan Institute of Biological Products/Sinopharm, Beijing Institute of Biological

524   Products/Sinopharm); five utilize the entire spike protein as the antigen

525   (Moderna/NIAID, CanSino Biological Inc./Beijing Institute of Biotechnology, University of

526   Oxford/AstraZeneca, Gamaleya Research Institute, Janssen Pharmaceutical

527   Companies) and one utilizes only the RBD (Pfizer/Fosun Pharma/BioNTech). In addition

528   to the greater sequence conservation of the spike protein relative to all other ORFs

529   outside of the polyprotein, it is the principal host-interacting protein of SARS-CoV-2,

530   making both the whole protein and the RBD obvious antigenic candidates. Most

531   mutations in the RBD were demonstrated to decrease infectivity, but some conferred

532   resistance to neutralizing antibodies(49). Multiple mutations in the RBD are signature

533   mutations in emerging variants and some have been demonstrated to result in

534   neutralizing antibody evasion(51). Different choices of the antigen could result in more

535   or less generalizable immunity to these variants.

536

## Conclusions

538   Virus evolution during a pandemic is a fast moving target, and unavoidably, aspects of

539   this analysis will be outdated by the time of publication. Nevertheless, several trends

540   revealed here appear general and robust. Although it is difficult to ascertain positive

541   selection for individual sites, the overall adaptive character of SARS-CoV-2 evolution

542   involving multiple amino acid replacements appears to be beyond reasonable doubt. As

543   expected, there are multiple positively selected sites in the S protein, but more

544   surprisingly, N protein includes several sites that appear to be strongly selected as well.

545   The involvement of these adaptive substitutions in the nuclear localization of the N

546 protein appears likely. Importantly, some of the mutations, for which positive selection

547 was inferred, co-occur on multiple occasions and seem to form a robust epistatic

548 network. Most likely, the effect of positive selection is manifested primarily at the level of

549 epistatic interactions.

550

551 Clearly, despite the dramatic reduction of global travel(76), the evolution of SARS-Cov-

552 2 is partly shaped by globalizing factors, including the increased virus fitness conferred

553 by S|D614G, N|R203K&G204R, and other positively selected substitutions. However,

554 we obtained strong evidence of both continuous virus diversification within geographic

555 regions and "speciation", that is, formation of stable, diverging region-specific variants.

556 This ongoing adaptive diversification could substantially prolong the pandemic and the

557 vaccination campaign, in which variant-specific vaccines are likely to be required.

558

559 **Author contributions**

560 EVK initiated the project; NDR and GF collected data; NDR, GF, YIW, PM, FZ and EVK
561 analyzed data; NDR and EVK wrote the manuscript that was edited and approved by all
562 authors.

20

## Acknowledgements

567

568

## Methods

570

### Multiple alignment of SARS-CoV-2 genomes

All available SARS-CoV-2 genomes as of January 8, 2021 were retrieved from the Genbank(13), Gisaid(14), and CNCB(15) datasets. Sequences with apparent anomalies (sequence inversion etc.) were immediately discarded. Sequences were harmonized to DNA (e.g. U was transformed to T to amend software compatibility) and clustered according to 100% identity with no coverage threshold using CD-HIT(77, 78), with ambiguous characters masked. All characters excepting ACGT were considered ambiguous. The least ambiguous sequence from each cluster was selected and sequences shorter than 25120 nucleotides were discarded.

Exterior ambiguous characters (preceding/succeeding the first/last defined nucleotide) were removed, and sequences with more than 10 remaining interior, ambiguous characters were discarded. A reference alignment was previously constructed using the same protocol as follows with the exception of the --keeplength specification in November, 2020. The updated database was aligned using multi-threaded MAFFT(79) with 80 cores (--thread 80, when more cores were allocated they were not utilized) and 3.8Tb of RAM to maintain usage of the normal DP algorithm(79) (--nomemsave) against this reference alignment (specifying --keeplength). Aligning "from scratch" without --keeplength proved to be prohibitively slow so we recommend first constructing a reference alignment from a suitable subset of sequences. Sequences sourced from non-human hosts were manually identified from the metadata and those excluded at the previous step were added to the alignment using MAFFT, (again specifying --keeplength). Note that use of the --keeplength option will not include insertions relative to the reference alignment.

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

The resulting alignment contained out-of-frame gaps. Gaps in the reference sequence, corresponding to insertions, were found to correspond to gaps in all but fewer than 1% of the remaining sequences (all gaps in the reference sequence correspond to gaps in the alignment from November, 2020, the use of --keeplength prohibited the recognition of any insertions relative to the reference sequence which were not present in this reference alignment). These sites were discarded. The remaining gaps, corresponding

605  to deletions relative to the reference sequence, shorter than three nucleotides were
606  replaced with the ambiguous character, N. Longer gaps were shifted into frame and
607  padded with ambiguous characters on either end of the gap, minimizing the number of
608  sites altered.

609  A fast, approximate tree was then built using FastTree(80) (parameters: -nt -gtr -gamma
610  -nosupport -fastest) to unambiguously define two clusters of sequences: an outgroup
611  consisting of 14 sequences sourced from non-human hosts prior to 2020 and the main
612  group. The tree construction requires the resolution of very short branch lengths which
613  makes it necessary to compile FastTree at double precision. Outliers from the remaining
614  sequences were then identified based on the Hamming distance (excluding gaps and
615  ambiguous characters) to the nearest neighbor, the Hamming distance to the
616  consensus, and the degree to which those substitutions relative to consensus were
617  clustered in the genome. At this step, 81 sequences were removed.

618  The resulting alignment, consisting of 98,185 sequences and 29,119 sites, was
619  maintained for the construction of the global tree and ancestral sequence
620  reconstruction. In an effort to minimize the impact of sequencing error on the tree
621  topology, as well as to decrease computational costs, a reduced alignment was then
622  constructed through the removal of 1) invariant sites, 2) sites invariant with the
623  exception of a single sequence, and 3) sites invariant throughout the main group with
624  the exception of at most one sequence representing each minority nucleotide.
625  Removing these sites created substantial redundancy, so a representative sequence
626  was selected for each cluster of 100% identity to yield an alignment consisting of 90,585
627  sequences and 16,487 sites. As described below and in the main text, a third alignment
628  was constructed including only the top 5% of sites with the most common substitutions
629  relative to consensus (of this second alignment) and again removing redundant
630  sequences to yield 32,563 sequences and 834 sites.

631

## Tree Construction

633  We sought to optimize tree topology with IQ-TREE(81); however, building the global
634  tree was computationally prohibitive, and thus, we proceeded to subsample the smallest
635  alignment (834 sites) as follows. First, a core set of maximally diverse sequences is
636  selected. The set is initialized with a pair of sequences: a sequence maximizing the
637  number of substitutions relative to consensus and a paired sequence which maximizes
638  the Hamming distance to itself. Sequences are then added to this core set one at a time
639  maximizing the minimum Hamming distance to any representative of the set until $N$
640  sequences are incorporated. Next, $ceil\left(L/(M-N)\right)$ resulting sets are initialized with this
641  core set where $M$ is the target number of sequences and $L$ is the total number of
642  sequences in the alignment (32,363). Then, sequences that have not yet been
643  incorporated into any resulting set are added to each resulting set, again one at a time,
644  maximizing the minimum distance to any representative of the set until $M$ sequences

23

645  are incorporated. The order of the resulting sets is randomized at each iteration without
646  repeats. Once every (main group) sequence has been incorporated into at least one
647  resulting set, sequences are randomly incorporated into each set until every set
648  contains *M* sequences. Finally, the outgroup is added to each resulting set. We chose
649  *M*=3,000 in an effort to optimize computational efficiency and *N*=300. Note that while
650  increasing *N* increases the number of sets required for alignment coverage, and thus
651  compute time, insufficient overlap between the sequences assigned each sub-alignment
652  greatly affects the results of subsequent steps. As discussed in the main text, executing
653  this protocol on an alignment containing most or all sites may not yield a consistent
654  deep tree topology or "skeleton" since maximizing the hamming distance of any subset
655  over all sites does not guarantee maximizing the tree distance in the resultant global
656  topology. This is why limiting the alignment to sites with common substitutions relative
657  to consensus is essential at this step.

658  A tree was then built, using IQ-TREE, for each maximally diverse set, with the
659  evolutionary model fixed to GTR+F+G4 and the minimum branch length decreased from
660  the default 10e-6 to 10e-7, according to the results of previous parameter studies(17).
661  These trees were then converted into constraint files and merged to generate a single
662  global constraint file for use within FastTree (parameters: -nt -gtr -gamma -cat 4 -
663  nosupport -constraints).

664  The remaining sequences excluded from this tree but present in the second alignment
665  (90,585 sequences and 16,487 sites) were then reintroduced as unresolved
666  multifurcations and a new constraint file from the multifurcated tree was constructed. A
667  second iteration of FastTree was initiated on the second alignment to produce an
668  intermediate tree. This tree was primarily constructed as an intermediate step to limit
669  the impact of sequencing errors on the final topology as mentioned in the main text;
670  however, it is also less computationally intensive. The last step was then repeated on
671  this intermediate tree to construct the global topology for the whole alignment. The final,
672  global tree was rooted at the outgroup.

673

## Reconstruction of Ancestral Genome Sequences

675  Ancestral states were estimated by Fitch Traceback(18). Briefly, character sets were
676  constructed from leaf to root where each node was assigned the intersection of the
677  descendant character sets if not empty and the union otherwise. Then, moving from root
678  to leaf, nodes with more than one character in their set were assigned the consensus
679  character if present in their set or a randomly chosen representative character
680  otherwise. Substitutions between states were identified and placed in the middle of the
681  branch bridging the pair of nodes.

682  Statistical associations between mutations were computed in a manner similar to that
683  previously described(35). Briefly, sequences were leaf-weighted based on the branch
684  lengths of the ultrameterized, tree. Every mutation present across the tree at 200 mean

24

685 leaf-weight equivalents or more was considered. The probability of independent co-
686 occurrence between any pair was estimated in two ways. An arbitrary member of the
687 pair was selected as the ancestral mutation, and the binomial probability:

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

688

689 was computed where $N\_total$ is the number of substitutions to the descendant mutation
690 across the entire ancestral record, $N\_pair$ is the number of substitutions to the
691 descendant which succeed or appear simultaneously with a substitution to the ancestral
692 mutation, and $F$ is the fraction of the tree (fraction of all applicable branch lengths)
693 occupied by the ancestral mutation. The ancestral/descendent designation was then
694 reversed and the "binomial score" was constructed as the negative log of the product of
695 these two terms. Additionally, for each pair, the observed and expected (product of the
696 tree fractions) tree intersections were calculated and the "Poisson score" (analogous to
697 the log-odds ratio) was calculated:

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

698 where PCDF(exp,obs) is the cumulative probability of a Poisson distribution with mean
699 "exp", the expected value of the data, and evaluated at "obs", the observed value of the
700 data. Both scores are reported. Table S3 displays putative positively selected mutations
701 with both scores above 5 or at least two simultaneous substitutions. Fig. 1D only
702 displays associations between mutations in the N or S proteins. Fig. S10 does not
703 exclude mutations with NCN context but meets all other statistical criteria for positive
704 selection and does not display mutations in the polyprotein.

705

706 ## Classical Multidimensional Scaling of the MSA

707 Pairwise Hamming distances were computed for all pairs of rows in the global MSA
708 ignoring gaps and ambiguous characters i.e. the sequences $X$="ATN-A" and
709 $Y$="NTAAT" would be assigned a distance of 1. The resulting distance matrix was
710 embedded in three dimensions with the MATLAB(82) routine "cmdscale". 100 rounds of
711 stochastically initiated k-means clustering of the embedding was conducted and the
712 optimum cluster number was determined to be 5 on the basis of the silhouette score
713 distribution (Fig S1).

714

715 ## Validation of Mutagenic Contexts

25

716   Mutations were divided into four categories: synonymous vs non-synonymous
717   substitutions and high vs low frequency of independent occurrence. For example,
718   consider codon X with 3 non-synonymous substitutions gat->ggt and 1 non-synonymous
719   substitution gat->cgt. In this context, a non-synonymous nucleotide substitution a->g of
720   frequency 4 would be recorded in nucleotide (X-1)*3+2. The low vs high frequency
721   threshold was determined by the $90^{th}$ percentile of the synonymous mutation frequency
722   distribution (operationally 7). For each mutation, the trinucleotide contexts from the
723   ancestral reconstruction at the nodes where the mutation occurred were compared to
724   the background genome-wide frequencies, computed for the inferred common ancestor
725   of SARS-CoV-2.

726

727   The expected frequencies of the trinucleotides using the background distribution were
728   tabulated; the Yates correction (+/-0.5 to the original count depending on whether the
729   count is below or above the expectation) was applied to the observed frequencies; the
730   log-odds ratios of the (corrected) observed frequencies to the expectation were
731   computed; and CMDS was applied to the Euclidean distances between the log-odds
732   vectors to embed the points onto a plane (Fig. S4 A.). This analysis was then repeated,
733   this time, distinguishing only between high and low frequency substitutions but not N
734   and S (Fig. S4 B). Finally, the differences in the contexts of high frequency synonymous
735   vs non-synonymous events were considered in the same manner and the chi-square
736   statistics ((observed-expected)^2/expected) were compared with the critical chi-square
737   value (p=0.05/64, df=1, Fig. S4 C.).

738

## Computation of *dN/dS*

740   For each of the 24 ORFs (splitting orf1ab into 15 segments corresponding to the 15
741   mature proteins, nsp11 and nsp12 combined), 10 reduced alignments were constructed
742   as follows. Sequences were ordered based on diversity, in the same order with which
743   they were included in the constraint trees. The first 10 sequences are conserved across
744   every alignment and the remaining 40 are unique to each alignment. The reference
745   sequence, NC_045512.2, was additionally added to each reduced alignment. PAML(83)
746   was then used to estimate tN, tS, *dN/dS*, N, S, and N/S for each segment and every
747   reduced alignment.

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

754

## Metadata Assignment

Headers for all isolates belonging to CD-HIT clusters with a representative incorporated into the alignment with fewer than 10 interior ambiguous characters were processed to extract the sequencing date and location. Sequencing location abbreviations were matched to full names and the latitude/longitude of a representative city for each location was retrieved from simplemaps (https://simplemaps.com/data/world-cities)(84).

## Regional Divergence Analysis

Two approaches, one partition dependent and one partition independent, were used as described in the main text. The Hellinger distance between regions over a sliding time window was computed between regions for the 11 (partitions/variant clades) group distribution. Next, 400 isolates were randomly selected from each region over a sliding window and 200 pairs within each region as well as 200 pairs between each pair of regions were composed. The tree distance between each pair was computed and the mean for each inter- and intra-regional pair tree-distance distribution was recorded. In Figs. 3C and S16, the $25^{th}$, $50^{th}$, and $75^{th}$ percentiles are shown of the 15 possible pairs of (6) regions. Regions are selected based on GISAID metadata. The inter-regional tree divergence (Figs. 3C, top and S16C) is reported as the ratio between the mean of the inter-regional pair tree-distance and the mean of the intra-regional pair tree distances across both regions.

**Figure legends**

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

**A.** Global tree reconstruction with 8 principal partitions and 3 variant clades enumerated and color-coded. **B.** Signatures of amino acid replacements for each partition. Sites are ordered as they appear in the genome. The proteins along with nucleotide and amino acid numbers are indicated underneath each column. **C.** Site history trees for spike 614 and nucleocapsid 203 positions. Nodes were included in this reduced tree based on the following criteria: those immediately succeeding a substitution; representing the last common ancestor of at least two substitutions; or terminal nodes representing branches of five sequences or more (approximately, based on tree weight). Edges are colored according to their position in the main partitions and the line type corresponds to the target mutation (solid) or any other state (dashed). Synonymous mutations are not shown. These sites are largely binary as are most sites in the genome. The sizes of the terminal node sizes are proportional to the log of the weight descendent from that node beyond which no substitutions in the site occurred. Node color corresponds to target mutation (black) or any other state (gray). **D.** Network of putative epistatic interactions for likely positively selected residues in the N and S proteins.

**Figure 2. Regional SARS-CoV-2 partition dynamics during the COVID-19 pandemic.** Probability distributions shown, for the absolute number of sequences, see Fig. S15.

**Figure 3. Global and regional trends in SARS-CoV-2 evolution. A.** Global distribution of sequences with sequencing locations in each of the six regions considered. Color scheme is for visual distinction only. **B.** Intra-regional diversity measured by the mean tree-distance for pairs of isolates. **C.** (Top) The Hellinger distance for all pairs of regions over the 11 partition/clade distribution. $25^{th}$, $50^{th}$, $75^{th}$ percentiles shown. (Bottom) The ratio of the mean tree-distance for pairs of isolates between regions vs. isolates within regions. $25^{th}$, $50^{th}$, $75^{th}$ percentiles shown. **D.** The frequency of S|614G, at least one NLS-associated variant (N|194L, N119L, N203K,

806    N205I, and N220V), and at least one emerging spike variant (Fig. S23, excluding

807    S|477N).

## Supplemental Figures

**Figure S1.** $25^{th}$, median (solid line), and $75^{th}$ percentiles of the silhouette score distribution for 100 stochastically initiated rounds of k-means clustering for 2-16 clusters and a projection of the 3D embedding of the pairwise Hamming distance matrix between SARC-CoV-2 genomes. Partitions are color-coded and wires enclose the convex hulls for each of the five optimal clusters.

**Figure S2. A.** Distributions of the moving average, respecting segment boundaries, across a 100 codon window for synonymous (blue) and amino acid (orange) substitutions. Solid lines: normal approximations of the distributions (same median and interquartile distance); solid lines: approximation with the same median and theoretical (Poisson) variance. **B.** Moving averages, respecting segment boundaries, across a 100 codon window for synonymous and nonsynonymous substitutions per site, raw (top) and normalized by the median (bottom). There are several regions in the genome with an apparent dramatic excess of synonymous substitutions: 5' end of orf1ab gene; most of the M gene; 3'-half of the N gene, as well as amino acid substitutions: most of the orf3a gene; most of the orf7a gene; most of the orf8 gene; and several regions in of the N gene.

**Figure S3.** Moving average over a window of 1000 codons, not respecting segment boundaries, of the total number of nucleotide exchanges n1->n2 summed over all substitutions. The ratio to the median over the entire alignment is also displayed as well as the normalized exchange distribution *(i.e.* #c->t/(#c->t+#c->g+#c->a)).

**Figure S4 A.** Two dimensional embedding of the Euclidean distances between the log-odds vectors of low and high frequency, nonsynonymous and synonymous mutations in the space of trinucleotide contexts relative to background expectation. The context of the high-frequency events (both S and N) is dramatically different from the background frequencies. There is a strong common component in the deviation of both kinds of high-frequency events. The context of the low-frequency events (both S and N) also differs slightly, in the same direction, from the background frequencies. There is a consistent distinction between synonymous and non-synonymous events, suggesting that a single mutagenic context or mechanistic bias does not account for both S and N events. **B.** Log odds ratio of low and high frequency mutations, both synonymous and nonsynonymous, relative to background expectation for each trinucleotide context. The NCN context (i.e. all mutations C->D) harbors dramatically more mutation events than the other contexts (all 16 NCN events are within the top 20 most-biased high-frequency events). The log-odds ratios for low-frequency events are poorly correlated with those

30

847    for high-frequency events, suggesting that different mechanisms may be responsible for
848    the strong bias observed among high frequency events and the weaker bias observed
849    among low frequency events. **C.** Log odds ratio of high frequency nonsynonymous
850    mutations relative to the background expectation from the sum of both high
851    synonymous and high nonsynonymous mutations vs. the sum + 1. There are 20
852    contexts where synonymous and non-synonymous events differ significantly (chi-sq>
853    11.28). 2/9 contexts with an excess of non-synonymous events are NCN (gct,tct). The
854    remaining 7 are NGN (agt,gga,aga,ggt,agc,tgt). This additionally suggests that these
855    non-synonymous events could be driven by other mechanisms. There is no correlation
856    between the frequency of event context and the log-odds ratio for non-synonymous
857    events, further suggesting that the log-odds ratio is not biased by hot-spot mutation
858    context.

859

860    **Figure S5.** Correspondence between the "tree length for dN", "tree length for dS", and
861    *dN/dS* between PAML and the results of the ancestral reconstruction utilizing Fitch
862    traceback across 24 ORFs. Three high outliers in the PAML tS distribution are identified
863    in the third plot and omitted from the first two.

864

865    **Figure S6. A.** The number of nonsynonymous events vs the number of synonymous
866    events per codon. **B.** The moving average of 100 codons, respecting segment
867    boundaries. **C.** The moving average after removing outlier high frequency events. Rho
868    refers to Spearman. Dashed lines are 2/1.3*x reflecting the genome-wide ratio of
869    nonsynonymous to synonymous substitutions, solid lines are linear best fit. Red points
870    correspond to the middle third of the N protein.

871

872    **Figure S7.** Moving averages across a 100 codon window for synonymous and
873    nonsynonymous substitutions per site in the N protein after removing outlier high
874    frequency events. The nonsynonymous substitution frequencies in the center of the
875    protein are not elevated relative to either terminus.

876

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

880

881    **Figure S9.** Site history trees for spike 69 as drawn in Fig. 1C.

882

883   **Figure S10.** Epistatic network for the tree including mutations with NCN context and
884   meeting all other criteria for positive selection. Mutations in the polyprotein are not
885   displayed.

886

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

889

890   **Figures S12-14.** Global distribution of sequences. Color represents the number of
891   sequences from that location and size represents the fraction of sequences from the
892   clade displayed. Partition indices are in the top left corner of each map.

893

894   **Figure S15.** Regional SARS-CoV-2 partition dynamics during the COVID-19 pandemic
895   (absolute number of sequences shown in contrast to Fig. 2).

896

897   **Figure S16.** The mean tree distance between pairs of isolates **A.** from different regions,
898   **B.** within the same region (averaged over both regions in each pair) and **C.** The ratio
899   over time (see Methods). 25$^{th}$, 50$^{th}$, and 75$^{th}$ percentiles of all 15 pairs of 6 regions. The
900   ratio reported is between the mean of the inter-regional pair tree-distance and the mean
901   of the intra-regional pair tree distances across both regions for each pair of regions.

902

903   **Figure S17.** Regional distributions of major partitions in the global topology March vs.
904   July and July vs. November.

905

906   **Figure S18.** The frequencies of NLS-associated mutations N|194L, N119L, N203K,
907   N205I, and N220V over time and across geographic regions along with S|614G for
908   reference.

909

910   **Figure S19.** The Kullback-Leibler divergence and sequence logo for the 15 most
911   divergent codons in sequences sourced after October 15, 2020 from Oceania in
912   partition 7 vs. all sequences from Oceania in partition 7.

913

914   **Figure S20.** The Kullback-Leibler divergence and sequence logo for the 15 most
915   divergent codons in sequences sourced after November 1, 2020 from Oceania in
916   partition 6 vs. all sequences from Oceania in partition 6.

917

918  **Figure S21.** The Kullback-Leibler divergence and sequence logo for the 15 most
919  divergent codons in sequences sourced after November 1, 2020 from Africa in partition
920  6 vs. all sequences from Africa in partition 6.

921

922  **Figures S22-24.** The frequencies of variant-associated mutations in the spike protein
923  over time and geographic regions.

924

925  **Figure S25.** Regional SARS-CoV-2 variant clade dynamics during the COVID-19
926  pandemic (log of absolute number of sequences shown).

927

928  ## Supplemental Tables

929

930  **Table S1.** The list of all mutations either in the top 100 most commonly observed or top
931  100 with the greatest number of parallel substitutions ordered as they appear in the
932  genome.

933

934  **Table S2.** List of sites most likely to be evolving under positive selection. For List 2 the
935  average tree fraction descendant from each candidate positively selected amino acid
936  replacement must be sufficiently large(see Methods).

937

938  **Table S3.** All epistatic interactions among states meeting the criteria outlined in the
939  main text for likely positive selection with binomial/Poisson scores greater than 5 or at
940  least 2 simultaneous substitutions. Each mutation must have a minimum weight of
941  approximately 200 leaves and each pair, 100 leaves. Each pair is arbitrarily ordered and
942  the numbers of simultaneous, descendant, and independent substitutions are tabulated.

943

944  **Table S4. List 1**. List of variant mutations and variant IDs sorted by the number of
945  variant ID's assigned to each mutation. **List 2**. List of all pairs of mutations associated
946  with a single variant ID (internal variant ID's excluded. **List 3**. List of putative epistatic
947  interactions between variant mutations and other states in the tree.

948

33

949 **Table S5.** The number of isolates (out of approximately 175k) observed to bear at least
950 one substitution relative to the reference sequence, NC_045512.2, within the regions
951 specified. These regions are commonly used within PCR assays for diagnostic testing.

952

953

954

955 1.   Drake JW & Holland JJ (1999) Mutation rates among RNA viruses. *Proceedings of the National*
956      *Academy of Sciences* 96(24):13910-13913.
957 2.   Sanjuán R (2012) From molecular genetics to phylodynamics: evolutionary relevance of
958      mutation rates across viruses. *PLoS Pathog* 8(5):e1002685.
959 3.   Simmonds P, Aiewsakun P, & Katzourakis A (2019) Prisoners of war - host adaptation and its
960      constraints on virus evolution. *Nat Rev Microbiol* 17(5):321-328.
961 4.   Elena SF & Sanjuán R (2005) Adaptive value of high mutation rates of RNA viruses: separating
962      causes from consequences. *J Virol* 79(18):11555-11558.
963 5.   Wertheim JO & Kosakovsky Pond SL (2011) Purifying selection can obscure the ancient age of
964      viral lineages. *Molecular biology and evolution* 28(12):3355-3365.
965 6.   Jenkins GM, Rambaut A, Pybus OG, & Holmes EC (2002) Rates of molecular evolution in RNA
966      viruses: a quantitative phylogenetic analysis. *Journal of molecular evolution* 54(2):156-165.
967 7.   Holmes EC (2003) Patterns of intra-and interhost nonsynonymous variation reveal strong
968      purifying selection in dengue virus. *Journal of virology* 77(20):11296-11298.
969 8.   Jerzak G, Bernard KA, Kramer LD, & Ebel GD (2005) Genetic variation in West Nile virus from
970      naturally infected mosquitoes and birds suggests quasispecies structure and strong purifying
971      selection. *The Journal of general virology* 86(Pt 8):2175.
972 9.   Wolf YI, Viboud C, Holmes EC, Koonin EV, & Lipman DJ (2006) Long intervals of stasis punctuated
973      by bursts of positive selection in the seasonal evolution of influenza A virus. *Biology direct*
974      1(1):1-19.
975 10.  Bush RM, Bender CA, Subbarao K, Cox NJ, & Fitch WM (1999) Predicting the evolution of human
976      influenza A. *Science* 286(5446):1921-1925.
977 11.  Bush RM, Fitch WM, Bender CA, & Cox NJ (1999) Positive selection on the H3 hemagglutinin
978      gene of human influenza virus A. *Molecular biology and evolution* 16(11):1457-1465.
979 12.  Koirala A, Joo YJ, Khatami A, Chiu C, & Britton PN (2020) Vaccines for COVID-19: The current
980      state of play. *Paediatric respiratory reviews* 35:43-49.
981 13.  Benson DA*, et al.* (2012) GenBank. *Nucleic acids research* 41(D1):D36-D42.
982 14.  Elbe S & Buckland-Merrett G (2017) Data, disease and diplomacy: GISAID's innovative
983      contribution to global health. *Global Challenges* 1(1):33-46.
984 15.  Zhao W-M*, et al.* (2020) The 2019 novel coronavirus resource. *Yi chuan= Hereditas* 42(2):212-
985      221.
986 16.  Lanfear R (A global phylogeny of SARS-CoV-2 from GISAID data, including sequences deposited
987      up to 31-July-2020. 2020. *Zenodo*.
988 17.  Morel B*, et al.* (2020) Phylogenetic analysis of SARS-CoV-2 data is difficult. *bioRxiv*.
989 18.  Fitch WM (1971) Toward defining the course of evolution: minimum change for a specific tree
990      topology. *Systematic Biology* 20(4):406-416.
991 19.  Kumar S*, et al.* (2020) An evolutionary portrait of the progenitor SARS-CoV-2 and its dominant
992      offshoots in COVID-19 pandemic. *bioRxiv*.

993   20.   Forster P, Forster L, Renfrew C, & Forster M (2020) Phylogenetic network analysis of SARS-CoV-2
994         genomes. *Proceedings of the National Academy of Sciences* 117(17):9241-9243.
995   21.   Fountain-Jones NM*, et al.* (2020) Emerging phylogenetic structure of the SARS-CoV-2 pandemic.
996         *Virus evolution* 6(2):veaa082.
997   22.   Mavian C*, et al.* (2020) Sampling bias and incorrect rooting make phylogenetic network tracing
998         of SARS-COV-2 infections unreliable. *Proceedings of the National Academy of Sciences*
999         117(23):12522-12523.
1000  23.   Sánchez-Pacheco SJ, Kong S, Pulido-Santacruz P, Murphy RW, & Kubatko L (2020) Median-
1001        joining network analysis of SARS-CoV-2 genomes is neither phylogenetic nor evolutionary.
1002        *Proceedings of the National Academy of Sciences* 117(23):12518-12519.
1003  24.   Pipes L, Wang H, Huelsenbeck J, & Nielsen R (2020) Assessing uncertainty in the rooting of the
1004        SARS-CoV-2 phylogeny. *bioRxiv*.
1005  25.   van Dorp L*, et al.* (2020) Emergence of genomic diversity and recurrent mutations in SARS-CoV-
1006        2. *Infection, Genetics and Evolution* 83:104351.
1007  26.   Timani KA*, et al.* (2005) Nuclear/nucleolar localization properties of C-terminal nucleocapsid
1008        protein of SARS coronavirus. *Virus research* 114(1-2):23-34.
1009  27.   Goldman N & Yang Z (1994) A codon-based model of nucleotide substitution for protein-coding
1010        DNA sequences. *Molecular biology and evolution* 11(5):725-736.
1011  28.   Richard D, Owen CJ, van Dorp L, & Balloux F (2020) No detectable signal for ongoing genetic
1012        recombination in SARS-CoV-2. *bioRxiv*:2020.2012.2015.422866.
1013  29.   Ye ZW*, et al.* (2020) Zoonotic origins of human coronaviruses. *Int J Biol Sci* 16(10):1686-1697.
1014  30.   Li X*, et al.* (2020) Emergence of SARS-CoV-2 through recombination and strong purifying
1015        selection. *Sci Adv* 6(27).
1016  31.   Zhu Z, Meng K, & Meng G (2020) Genomic recombination events may reveal the evolution of
1017        coronavirus and the origin of SARS-CoV-2. *Sci Rep* 10(1):21617.
1018  32.   Ji W, Wang W, Zhao X, Zai J, & Li X (2020) Cross-species transmission of the newly identified
1019        coronavirus 2019-nCoV. *J Med Virol* 92(4):433-440.
1020  33.   Graham RL & Baric RS (2010) Recombination, reservoirs, and the modular spike: mechanisms of
1021        coronavirus cross-species transmission. *J Virol* 84(7):3134-3146.
1022  34.   Bobay LM, O'Donnell AC, & Ochman H (2020) Recombination events are concentrated in the
1023        spike protein region of Betacoronaviruses. *PLoS Genet* 16(12):e1009272.
1024  35.   Rochman ND, Wolf YI, & Koonin EV (2020) Deep phylogeny of cancer drivers and compensatory
1025        mutations. *Communications biology* 3(1):1-11.
1026  36.   Korber B*, et al.* (2020) Tracking changes in SARS-CoV-2 Spike: evidence that D614G increases
1027        infectivity of the COVID-19 virus. *Cell* 182(4):812-827. e819.
1028  37.   Volz E*, et al.* (2021) Evaluating the effects of SARS-CoV-2 Spike mutation D614G on
1029        transmissibility and pathogenicity. *Cell* 184(1):64-75. e11.
1030  38.   Breen MS, Kemena C, Vlasov PK, Notredame C, & Kondrashov FA (2012) Epistasis as the primary
1031        factor in molecular evolution. *Nature* 490(7421):535-538.
1032  39.   Starr TN & Thornton JW (2016) Epistasis in protein evolution. *Protein Sci* 25(7):1204-1218.
1033  40.   Phillips PC (2008) Epistasis--the essential role of gene interactions in the structure and evolution
1034        of genetic systems. *Nat Rev Genet* 9(11):855-867.
1035  41.   Domingo J, Baeza-Centurion P, & Lehner B (2019) The Causes and Consequences of Genetic
1036        Interactions (Epistasis). *Annu Rev Genomics Hum Genet* 20:433-460.
1037  42.   Faria NR*, et al.* (2021) Genomic characterisation of an emergent SARS-CoV-2 lineage in Manaus:
1038        preliminary findings. *January* 12:2021.

1039 43. Tegally H*, et al.* (2020) Emergence and rapid spread of a new severe acute respiratory
1040     syndrome-related coronavirus 2 (SARS-CoV-2) lineage with multiple spike mutations in South
1041     Africa. *medRxiv*.
1042 44. Voloch CM*, et al.* (2020) Genomic characterization of a novel SARS-CoV-2 lineage from Rio de
1043     Janeiro, Brazil. *medRxiv*.
1044 45. Zhang W*, et al.* (2021) Emergence of a Novel SARS-CoV-2 Variant in Southern California. *JAMA*.
1045 46. England PH (Variants: distribution of cases data.
1046 47. Brejová B*, et al.* (2021) B. 1.258_Delta, a SARS-CoV-2 variant with Delta_H69, Delta_V70 in the
1047     Spike protein circulating in the Czech Republic and Slovakia. *arXiv preprint arXiv:2102.04689*.
1048 48. Statens Serum Institut (2020) Mutations in the mink virus.
1049     https://www.ssi.dk/aktuelt/nyheder/2020/mutationer-i-minkvirus
1050 49. Munnink BBO*, et al.* (2021) Transmission of SARS-CoV-2 on mink farms between humans and
1051     mink and back to humans. *Science* 371(6525):172-177.
1052 50. Li Q*, et al.* (2020) The impact of mutations in SARS-CoV-2 spike on viral infectivity and
1053     antigenicity. *Cell* 182(5):1284-1294. e1289.
1054 51. Wibmer CK*, et al.* (2021) SARS-CoV-2 501Y. V2 escapes neutralization by South African COVID-19
1055     donor plasma. *BioRxiv*.
1056 52. Fiorentini S*, et al.* (2021) First detection of SARS-CoV-2 spike protein N501 mutation in Italy in
1057     August, 2020. *The Lancet. Infectious Diseases*.
1058 53. van Dorp L*, et al.* (2020) Recurrent mutations in SARS-CoV-2 genomes isolated from mink point
1059     to rapid host-adaptation. *bioRxiv*.
1060 54. Thomson EC*, et al.* (2020) The circulating SARS-CoV-2 spike variant N439K maintains fitness
1061     while evading antibody-mediated immunity. *bioRxiv*.
1062 55. Rochman ND, Wolf YI, & Koonin EV (2020) Substantial Impact of Post Vaccination Contacts on
1063     Cumulative Infections during Viral Epidemics. *medRxiv*.
1064 56. Gandon S & Day T (2007) The evolutionary epidemiology of vaccination. *Journal of the Royal
1065     Society Interface* 4(16):803-817.
1066 57. Brueggemann AB, Pai R, Crook DW, & Beall B (2007) Vaccine escape recombinants emerge after
1067     pneumococcal vaccination in the United States. *PLoS Pathog* 3(11):e168.
1068 58. Scherer A & McLean A (2002) Mathematical models of vaccination. *British Medical Bulletin*
1069     62(1):187-199.
1070 59. Geoghegan JL & Holmes EC (2018) The phylogenomics of evolving virus virulence. *Nature
1071     Reviews Genetics* 19(12):756-769.
1072 60. Cressler CE, McLeod DV, Rozins C, Van Den Hoogen J, & Day T (2016) The adaptive evolution of
1073     virulence: a review of theoretical predictions and empirical tests. *Parasitology* 143(7):915-930.
1074 61. Rochman ND, Wolf YI, & Koonin EV (2020) Evolution of Human Respiratory Virus Epidemics.
1075     *medRxiv*.
1076 62. Gussow AB*, et al.* (2020) Genomic determinants of pathogenicity in SARS-CoV-2 and other
1077     human coronaviruses. *Proceedings of the National Academy of Sciences* 117(26):15193-15199.
1078 63. Hodcroft EB*, et al.* (2020) Emergence and spread of a SARS-CoV-2 variant through Europe in the
1079     summer of 2020. *MedRxiv*.
1080 64. Zhang Y*, et al.* (2020) The ORF8 protein of SARS-CoV-2 mediates immune evasion through
1081     potently downregulating MHC-I. *biorxiv*.
1082 65. Zinzula L (2020) Lost in deletion: The enigmatic ORF8 protein of SARS-CoV-2. *Biochemical and
1083     Biophysical Research Communications*.
1084 66. Gong LI, Suchard MA, & Bloom JD (2013) Stability-mediated epistasis constrains the evolution of
1085     an influenza protein. *Elife* 2:e00631.

1086   67.   Sanjuán R, Cuevas JM, Moya A, & Elena SF (2005) Epistasis and the adaptability of an RNA virus.
1087         *Genetics* 170(3):1001-1008.
1088   68.   Lyons DM & Lauring AS (2018) Mutation and Epistasis in Influenza Virus Evolution. *Viruses* 10(8).
1089   69.   Kryazhimskiy S, Dushoff J, Bazykin GA, & Plotkin JB (2011) Prevalence of epistasis in the
1090         evolution of influenza A surface proteins. *PLoS Genet* 7(2):e1001301.
1091   70.   Kuipers J*, et al.* (2020) Within-patient genetic diversity of SARS-CoV-2. *BioRxiv*.
1092   71.   Rose R*, et al.* (2020) Intra-host site-specific polymorphisms of SARS-CoV-2 is consistent across
1093         multiple samples and methodologies. *medRxiv*:2020.2004.2024.20078691.
1094   72.   Armero A, Berthet N, & Avarre JC (2021) Intra-Host Diversity of SARS-Cov-2 Should Not Be
1095         Neglected: Case of the State of Victoria, Australia. *Viruses* 13(1).
1096   73.   Artesi M*, et al.* (2020) A recurrent mutation at position 26340 of SARS-CoV-2 is associated with
1097         failure of the E gene quantitative reverse transcription-PCR utilized in a commercial dual-target
1098         diagnostic assay. *Journal of clinical microbiology* 58(10).
1099   74.   Ortiz-Prado E, et al. (2020) Clinical, molecular and epidemiological characterization of the SARS-
1100         CoV2 virus and the Coronavirus disease 2019 (COVID-19), a comprehensive literature review.
1101         *Diagnostic microbiology and infectious disease* 115094.
1102   75.   Dong Y, et al. (2020) A systematic review of SARS-CoV-2 vaccine candidates. *Signal transduction*
1103         *and targeted therapy* 5.1.
1104   76.   Lai S*, et al.* (2020) Assessing the effect of global travel and contact reductions to mitigate the
1105         COVID-19 pandemic and resurgence. *medRxiv*.
1106   77.   Fu L, Niu B, Zhu Z, Wu S, & Li W (2012) CD-HIT: accelerated for clustering the next-generation
1107         sequencing data. *Bioinformatics* 28(23):3150-3152.
1108   78.   Li W & Godzik A (2006) Cd-hit: a fast program for clustering and comparing large sets of protein
1109         or nucleotide sequences. *Bioinformatics* 22(13):1658-1659.
1110   79.   Katoh K, Misawa K, Kuma Ki, & Miyata T (2002) MAFFT: a novel method for rapid multiple
1111         sequence alignment based on fast Fourier transform. *Nucleic acids research* 30(14):3059-3066.
1112   80.   Price MN, Dehal PS, & Arkin AP (2010) FastTree 2–approximately maximum-likelihood trees for
1113         large alignments. *PLoS one* 5(3):e9490.
1114   81.   Nguyen L-T, Schmidt HA, Von Haeseler A, & Minh BQ (2015) IQ-TREE: a fast and effective
1115         stochastic algorithm for estimating maximum-likelihood phylogenies. *Molecular biology and*
1116         *evolution* 32(1):268-274.
1117   82.   MathWorks I (1992) *MATLAB, high-performance numeric computation and visualization*
1118         *software: reference guide* (MathWorks).
1119   83.   Yang Z (2007) PAML 4: phylogenetic analysis by maximum likelihood. *Molecular biology and*
1120         *evolution* 24(8):1586-1591.
1121   84.   Simplemaps. World Cities Database. https://simplemaps.com/data/world-cities

1122

**A** North America
**B** Europe
**C** Asia
**D** South America
**E** Africa
**F** Oceania

partitions

| | |
|---|---|
| 1 | 7 |
| 2 | 8 |
| 3 | v1 |
| 4 | v2 |
| 5 | v3 |
| 6 | |