

1 The global population of SARS-CoV-2 is composed of six major subtypes

2

3 Ivair José Morais Júnior^{ad}, Richard Costa Polveiro^b, Gabriel Medeiros Souza^a, Daniel Inserra
4 Bortolin^a, Flávio Tetsuo Sasaki^c, Alison Talis Martins Lima^{a#}

5 ^aInstituto de Ciências Agrárias/Universidade Federal de Uberlândia, Uberlândia, MG 38410-337,
6 Brazil

7 ^bDepartamento de Veterinária/Universidade Federal de Viçosa, Viçosa, MG 36570-900, Brazil

8 ^cInstituto de Biotecnologia/Universidade Federal de Uberlândia, Monte Carmelo, MG 38500-000,
9 Brazil

10

11 ^dCurrent address: Departamento de Fitopatologia/Universidade de Brasília, Brasília, DF 73345-
12 010, Brazil

13

14 #Corresponding author: Alison Talis Martins Lima

15 Phone: (+55-34) 2512-6716; E-mail: atmlima@ufu.br

16

17 Abstract

18

19 The World Health Organization characterized the COVID-19 as a pandemic in March 2020, the second
20 pandemic of the 21st century. Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) is a
21 positive-stranded RNA betacoronavirus of the family *Coronaviridae*. Expanding virus populations, as
22 that of SARS-CoV-2, accumulate a number of narrowly shared polymorphisms imposing a
23 confounding effect on traditional clustering methods. In this context, approaches that reduce the
24 complexity of the sequence space occupied by the SARS-CoV-2 population are necessary for a robust
25 clustering. Here, we proposed the subdivision of the global SARS-CoV-2 population into sixteen well-
26 defined subtypes by focusing on the widely shared polymorphisms in nonstructural (*nsp3*, *nsp4*, *nsp6*,
27 *nsp12*, *nsp13* and *nsp14*) cistrons, structural (*spike* and *nucleocapsid*) and accessory (*ORF8*) genes.
28 Six virus subtypes were predominant in the population, but all sixteen showed amino acid
29 replacements which might have phenotypic implications. We hypothesize that the virus subtypes
30 detected in this study are records of the early stages of the SARS-CoV-2 diversification that were
31 randomly sampled to compose the virus populations around the world, a typical founder effect. The
32 genetic structure determined for the SARS-CoV-2 population provides substantial guidelines for
33 maximizing the effectiveness of trials for testing the candidate vaccines or drugs.

34 Main

35 In December 2019, a local pneumonia outbreak of initially unknown etiology was detected in
36 Wuhan (Hubei, China) and quickly determined to be caused by a novel coronavirus¹, named severe
37 acute respiratory syndrome coronavirus 2 (SARS-CoV-2)² and the disease as COVID-19³. SARS-
38 CoV-2 is classified in the family *Coronaviridae*, genus *Betacoronavirus*, which comprises enveloped,
39 positive stranded RNA viruses of vertebrates². Two-thirds of SARS-CoVs genome is covered by the
40 ORF1ab, that encodes a large polypeptide which is cleaved into 16 nonstructural proteins (NSPs)
41 involved in replication-transcription in vesicles from endoplasmic reticulum (ER)-derived
42 membranes^{4,5}. The last third of the virus genome encodes four essential structural proteins: spike (S),
43 envelope (E), membrane (M), nucleocapsid (N) and several accessory proteins that interfere with the
44 host innate immune response⁶.

45 Populations of RNA viruses evolve rapidly due to their large population sizes, short generation
46 times, and high mutation rates of viruses, this latter is a consequence of the RNA-dependent RNA
47 polymerase (RdRP) which lacks the proofreading activity⁷. In fact, virus populations are composed of
48 a broad spectrum of closely related genetic variants resembling one or more master sequences⁸⁻¹⁰.
49 Mutation rates inferred for SARS-CoVs are considered moderate^{11,12} due to the independent
50 proofreading activity¹³. However, the large SARS-CoV genomes (from 27 to 31 kb)¹⁴ provide to them
51 the ability to explore the sequence space¹⁵. In order to better understand the diversification of SARS-
52 CoV-2 genomes during the pandemics (from December 2019 to March 25, 2020), we applied a simple,
53 but robust approach to reduce the complexity of the sequence space occupied by the virus population
54 by detecting its widely shared polymorphisms.

55 The 767 SARS-CoV-2 genomes with high sequencing coverage obtained from GISAID
56 (<https://www.gisaid.org/>) and GenBank were clustered into 593 haplotypes (Supplementary Table 1).
57 We conducted a fine-scale sequence variation analysis on the 593 genomes-containing alignment by
58 calculating the nucleotide diversity (π) using a sliding window and step size of 300 and 20 nucleotides,
59 respectively (multiple sequence alignments generated in this study are available from the authors upon
60 request). Such an approach allows to identify genomic regions of increased genetic variation
61 prevenient from polymorphic sites harboring two or more distinct nucleotide bases. Noticeably, one
62 or more large clusters of closely related sequences, when analyzed by this approach, show locally
63 increased nucleotide diversity. We observed a contrasting distribution of the genetic variation across
64 the full-length genomes of SARS-CoV-2 (Figure 1) with eight segments showing increased genetic

65 variation, arbitrarily defined as nucleotide (nt) segments with $\pi \geq 0.001$ (Table 1). Seven out of eight
66 segments had about 280 nucleotides in length, corresponding approximately to the size of a single
67 sliding window, except the S10 whose length was equivalent to two sliding windows (600 nt). To
68 further investigate the diversification of segments with contrasting content of genetic variability, we
69 constructed maximum likelihood (ML) phylogenetics trees and analyzed the diversification patterns
70 of eight segments with higher (S2, 4, 6, 8, 10, 12, 14 and 16), and nine with lower (S1, 3, 5, 7, 9, 11,
71 13, 15 and 17) content of genetic variation, respectively (the ML analyses were used in this study
72 essentially as a clustering method due to the weak phylogenetic signal in the data set).

73 Although the data set was composed of hundreds of SARS-CoV-2 genomes sampled from
74 around the world, in the S2-based tree we observed two clusters (Figure 2). Markedly, each cluster
75 was composed of very closely related, if not identical, sequences. Therefore, the increased content of
76 genetic variation at the S2 was a result of the inter-cluster sequence comparisons. Similar results were
77 obtained for the other seven ML-trees based on segments with increased genetic variation
78 (Supplementary Figure 1). In contrast, the ML-trees based on segments with lower content of genetic
79 variation did not show a consistent number of well-defined clusters (Supplementary Figure 2).

80 We mapped the polymorphic sites in segments whose trees showed two well-defined clusters
81 (Table 1). Only a few (from one to three) nt positions with polymorphisms shared by a number of
82 SARS-CoV-2 genomes could be identified within each segment with increased genetic variation.
83 These polymorphisms were henceforth referred to as ‘widely shared polymorphisms’ (WSPs), while
84 the remaining nt sites in virus genomes were designed as ‘non widely shared polymorphisms’
85 (nWSPs).

86 We compared the topologies of the seventeen ML-trees (eight based on WSPs-containing
87 segments and nine based on segments with nWSPs, Supplementary Figures 1 and 2, respectively). The
88 topology of trees based on WSPs-containing segments was considerably congruent indicating a
89 frequent, but not strict, co-segregation of nt at WSPs. For example, the SARS-CoV-2 reference isolate
90 (GISAID accession ID: EPI_ISL_402124, GenBank accession: MN908947) was placed in an
91 equivalent major cluster in all WSPs-containing segments-based ML-trees. In contrast, the isolate
92 EPI_ISL_416036 was placed in an equivalent minor cluster for S2, S8, S12 and S16-based ML-trees
93 and in the major cluster for S4, S6, S10 and S14-based ML-trees. On the contrary, the topology of
94 those ML-trees based on segments with nWSPs were highly incongruent suggesting that such regions
95 represent a wide mutant spectrum of narrowly shared polymorphisms. It is important to note that there

96 are minor clusters in nWSPs-containing segments-based ML-trees, *e.g.*, in those for S1, S13 and S17.
97 This is a consequence of our conservative threshold in which we focused on segments with $\pi \geq 0.001$.
98 S1, S13 and S17 also show locally increased genetic variation with π higher than 0.0005 but lower
99 than 0.001. For example, stretches 889 - 1,169; 1,409 - 1,509 (within the S1); 25,403 - 25,693 (S13);
100 29,538 - 29,610 (S17).

101 Therefore, our approach reduced the complexity of the sequence space occupied by the SARS-
102 CoV-2 genomes and provided a robust clustering solution based on the combination of 12 WSPs
103 (Table 1) to identify the major viral genotypes spread worldwide (Figure 3). The global population of
104 SARS-CoV-2 is structured into six major subtypes (I - VI), comprising 578 out of 593 (about 97.5%)
105 isolates analyzed in this study. The Subtype I (N=132) was represented by the combination of the most
106 frequent nucleotides at all WSPs, *i.g.*, the canonical genotype: CCGCCACAUGGG. The SARS-CoV-
107 2 reference isolate is a representative member of this subtype. Subtype IV (N=91) was represented by
108 the combination of the most frequent nucleotides at eleven out of 12 WSPs (**CCUCCACAUGGG**;
109 the most frequent nucleotides at each WSP are highlighted in bold and underlined). Subtypes V (N=74,
110 **CUGCCACACGGG**), II (N=122, **UCGUCACGUGGG**), III (N=101, **CUGCUGUAACGGG**) and VI
111 (N=58, **UCGUCACGUAAAC**) were represented by the combination of the most frequent nucleotides
112 at ten, nine, seven and six out of 12 WSPs, respectively. It is important to emphasize that the
113 contrasting sample sizes (Subtypes I - IV *vs.* VII - XVI) are not necessarily associated with fitness
114 variation and might be due to a sampling bias. For example, the three isolates composing the Subtype
115 VIII showed a genotype (**CUGCCAUAACGGG**) very similar to that of the canonical reference isolate.
116 In addition, even though our data set was composed exclusively by genomes with high sequencing
117 coverage, we cannot rule out that the virus subtypes X to XVI, which were represented by a single
118 genome might be a consequence of poor sampling or sequencing errors.

119 The phylogenetic tree depicting all 593 SARS-CoV-2 haplotypes showed some geographical
120 structure with two clusters: a smaller one comprised of isolates mostly sampled from Western
121 hemisphere (Subtypes II, VI, IX, X and XI) and a larger one whose isolates were sampled from
122 Western and Eastern hemispheres (I, III, IV, V, VII, VIII, XII, XIII, XIV, XV and XVI).

123 We hypothesize that our clustering method for the SARS-CoV-2 population could involve at
124 some extent a biological context. Nine out of 12 WSPs led to amino acid replacements (Table 2), *e.g.*,
125 the WSP *nsp6*-[111] in nine SARS-CoV-2 subtypes led to a leucine at the aa residue#37 of the protein
126 and a phenylalanine in seven other subtypes. NSP6 is an integral membrane protein that interferes in

127 the autophagosome formation during the SARS-CoV infection. Additionally, in yeast two-hybrid
128 experiments, NSP6 has been shown to interact with NSP3¹⁶. Some evidence demonstrates that NSP6
129 protein limits the expansion of autophagosomes or, alternatively, might remove host proteins involved
130 in inhibition of viral replication by activating autophagy from the ER¹⁷.

131 The WSP *nsp12*-[967] resulted in a proline in eleven subtypes of SARS-CoV-2 and a leucine
132 in others five subtypes at the aa residue #323 of the NSP12 (RNA-dependent RNA polymerase, RdRP)
133 protein. It is located at the Interface domain of RdRP of SARS-CoV-2, which is responsible for the
134 connection between the nidovirus RdRP-associated nucleotidyltransferase domain (NiRAN) and the
135 “Right hand” polymerase domain¹⁸. The S protein mediates viral entry into host cells by first binding
136 to a receptor, angiotensin-converting enzyme 2 (ACE2), through the receptor-binding domain (RBD)
137 in the S1 subunit and then fusing the viral and host membranes through the S2 subunit^{19–22}. Sites of
138 glycosylation are important for S protein folding²³, affecting priming by host proteases²⁴ and might
139 modulate antibody recognition^{25,26}. The WSP *S*-[1,841] resulted in a glycine and an aspartate at the aa
140 residue#614 of the S protein in six and ten subtypes of SAR-CoV-2, respectively. The replacement
141 was mapped in the intermediate region between the S1 and S2 subunits. This WSP is near a
142 glycosylation site (N616CT)²⁷.

143 The WSP *ORF8*-[251] involved a non-synonymous mutation at the codon#84 encoding for
144 leucine and serine in nine and seven subtypes, respectively. The SARS-CoV ORF8 encodes for an
145 ER-associated protein that induces the activation of ATF6, and this latter is an ER stress-regulated
146 transcription factor that stimulates the production of chaperones²⁸. In addition, the ORF8 protein has
147 been demonstrated to induce apoptosis²⁹. In SARS epidemics, the *ORF8* from different coronaviruses
148 was targeted by a number of mutations and recombination events during transmission from animals to
149 humans³⁰.

150 Three WSPs mapped in the *N* gene led to two amino acid replacements at residues#203 and
151 #204. The multifunctional N protein is composed of three domains³¹, two of which are structurally
152 independent: the N-terminal domain (NTD) and the C-terminal domain (CTD). Both amino acid
153 replacements were mapped in an intermediary domain referred to as the linker region (LKR), a
154 positively charged serine-arginine-rich region. As an intrinsically disordered region (IDR) it allows
155 the independent folding of the NTD and CTD³² and is also functionally implicated in RNA binding
156 activity³¹. Key determinants of the interaction between the N and NSP3 proteins were also mapped at
157 the LKR³³. The SARS-CoV N protein is also responsible for an antigenic response in humans

158 predominantly involving the immunoglobulin G³⁴ (the three-dimensional protein maps highlighting
159 the amino acid replacements are available as Supplementary Figure 3). Although the host biological
160 factors involved in the response to SARS-CoV-2 infection are still poorly known, the existence of
161 distinct virus subtypes, all of them exhibiting amino acid replacements, could alter important aspects
162 of COVID-19.

163 We hypothesized that in the early stages of the SARS-CoV-2 epidemics, due to the rapid virus
164 population expansion, a number of genetic variants might have arisen followed by a spread of non-
165 representative sampling of variants to other countries and continents, *i.g.*, a founder effect. We argue
166 that the virus subtypes and their associated WSPs detected in this study would be records of
167 diversification in these early stages of the epidemics after transmission from animal to humans. After
168 the virus introduction in a given geographic region, a number of unique or narrowly shared mutations
169 is accumulated, however, most of them reduce the fitness and are removed by purifying selection in a
170 medium to long term evolutionary scale, tending to a decreasing genetic variability⁸.

171 We propose a classification into at least sixteen distinct subtypes of SARS-CoV-2, six of them
172 accounting for more than 97% of the sampled isolates from around the world. Such classification
173 might guide the validation of candidate vaccines or drugs for the widest range of virus subtypes. In
174 this context, our clustering solution provides a robust approach to effectively reduce the complexity
175 of the mutant spectrum composed of closely related SARS-CoV-2 genomes focusing on WSPs.
176 Additionally, through the exhaustive sequencing, it would be possible to identify novel virus subtypes
177 and follow the evolutionary dynamics of the SARS-CoV-2 population during the adaptive process
178 imposed by the human host.

179 **Table 1.** Characterization of the WSPs detected in genomes of SARS-CoV-2
 180

Segment ID	Segment position* (begin - end)	WSPs†	nt mutation (# isolates)	Position in the codon	#codon	Amino acid
S2	2,872 - 3,152	<i>nsp3</i> -[318]	U (184) / C (409)	Third	106	Phenylalanine/Phenylalanine
S4	8,612 - 8,892	<i>nsp4</i> -[228]	U (183) / C (410)	Third	76	Serine/Serine
S6	10,932 - 11,192	<i>nsp6</i> -[111]	C (1) / U (99) / G (493)	Third	37	Phenylalanine/Phenylalanine/Leucine
S8	14,232 - 14,512	<i>nsp12</i> -[967]	U (184) / C (409)	Second	323	Leucine/Proline
S10	17,573 - 18,173	<i>nsp13</i> -[1,511]	U (101) / C (492)	Second	504	Leucine/Proline
		<i>nsp13</i> -[1,622]	G (101) / A (492)	Second	541	Cysteine/Tyrosine
		<i>nsp14</i> -[21]	U (105) / C (488)	Third	7	Leucine/Leucine
S12	23,243 - 23,523	<i>S</i> -[1,841]	G (185) / A (408)	Second	614	Glycine/Aspartate
S14	27,977 - 28,258	<i>ORF8</i> -[251]	C (184) / U (409)	Second	84	Serine/Leucine
S16	28,718 - 28,998	<i>N</i> -[608]	A (60) / G (533)	Second	203	Lysine/Arginine
		<i>N</i> -[609]	A (60) / G (533)	Third		
		<i>N</i> -[610]	C (60) / G (533)	First	204	Glycine/Arginine

181 *Relative to the multiple sequence alignment constructed for full-length genomes

182 †Widely Shared Polymorphisms (WSPs) are conventionally indicated by *cistron/gene*-[nt position within the cistron or gene]

183 **Table 2.** Amino acid composition of each virus subtype at WSP positions

184

SUBTYPE	N*	WSP POSITION											
		<i>nsp3</i>	<i>nsp4</i>	<i>nsp6</i>	<i>nsp12</i>	<i>nsp13</i>		<i>nsp14</i>	<i>S</i>	<i>ORF8</i>	<i>N</i>		
		#318	#228	#111	#967	#1511	#1622	#21	#1,841	#251	#608	#609	#610
I	132	C [Phe] [†]	C [Ser]	G [Leu]	C [Pro]	C [Pro]	A [Tyr]	C [Leu]	A [Asp]	U [Leu]	G [Arg]	G [Arg]	G [Arg]
II	122	U [Phe]	C [Ser]	G [Leu]	U [Leu]	C [Pro]	A [Tyr]	C [Leu]	G [Gly]	U [Leu]	G [Arg]	G [Arg]	G [Arg]
III	101	C [Phe]	U [Ser]	G [Leu]	C [Pro]	U [Leu]	G [Cys]	U [Leu]	A [Asp]	C [Ser]	G [Arg]	G [Arg]	G [Arg]
IV	91	C [Phe]	C [Ser]	U [Phe]	C [Pro]	C [Pro]	A [Tyr]	C [Leu]	A [Asp]	U [Leu]	G [Arg]	G [Arg]	G [Arg]
V	74	C [Phe]	U [Ser]	G [Leu]	C [Pro]	C [Pro]	A [Tyr]	C [Leu]	A [Asp]	C [Ser]	G [Arg]	G [Arg]	G [Arg]
VI	58	U [Phe]	C [Ser]	G [Leu]	U [Leu]	C [Pro]	A [Tyr]	C [Leu]	G [Gly]	U [Leu]	A [Lys]	A [Lys]	C [Gly]
VII	3	C [Phe]	U [Ser]	U [Phe]	C [Pro]	C [Pro]	A [Tyr]	C [Leu]	A [Asp]	C [Ser]	G [Arg]	G [Arg]	G [Arg]
VIII	3	C [Phe]	U [Ser]	G [Leu]	C [Pro]	C [Pro]	A [Tyr]	U [Leu]	A [Asp]	C [Ser]	G [Arg]	G [Arg]	G [Arg]
IX	2	U [Phe]	C [Ser]	U [Phe]	U [Leu]	C [Pro]	A [Tyr]	C [Leu]	G [Gly]	U [Leu]	A [Lys]	A [Lys]	C [Gly]
X	1	U [Phe]	C [Ser]	U [Phe]	U [Leu]	C [Pro]	A [Tyr]	C [Leu]	G [Gly]	U [Leu]	G [Arg]	G [Arg]	G [Arg]
XI	1	U [Phe]	C [Ser]	G [Leu]	C [Pro]	C [Pro]	A [Tyr]	C [Leu]	G [Gly]	U [Leu]	G [Arg]	G [Arg]	G [Arg]
XII	1	C [Phe]	U [Ser]	C [Phe]	C [Pro]	C [Pro]	A [Tyr]	C [Leu]	A [Asp]	C [Ser]	G [Arg]	G [Arg]	G [Arg]
XIII	1	C [Phe]	C [Ser]	G [Leu]	C [Pro]	C [Pro]	A [Tyr]	C [Leu]	A [Asp]	C [Ser]	G [Arg]	G [Arg]	G [Arg]
XIV	1	C [Phe]	U [Ser]	G [Leu]	C [Pro]	C [Pro]	A [Tyr]	C [Leu]	G [Gly]	C [Ser]	G [Arg]	G [Arg]	G [Arg]
XV	1	C [Phe]	C [Ser]	U [Phe]	U [Leu]	C [Pro]	A [Tyr]	C [Leu]	A [Asp]	U [Leu]	G [Arg]	G [Arg]	G [Arg]
XVI	1	C [Phe]	C [Ser]	U [Phe]	C [Pro]	C [Pro]	A [Tyr]	U [Leu]	A [Asp]	U [Leu]	G [Arg]	G [Arg]	G [Arg]

185 *Sample size

186 #Nucleotide position

187 †Nucleotide base and the encoded amino acid residue

188 **Figure legends**

189 **Figure 1.** Mean pairwise number of nucleotide differences per site (nucleotide diversity, π)
190 calculated using a sliding window of 300 nucleotides across the multiple sequence alignment for
191 full-length genomes of SARS-CoV-2. The red dashed line at $\pi = 0.001$ represents an arbitrary
192 threshold to subdivide the segments (S) with increased (S2, 4, 6, 8, 10, 12, 14 and 16) and lower
193 (S1, 3, 5, 7, 9, 11, 13,15 and 17) content of genetic variation. The SARS-CoV-2 genome
194 organization is represented on top of the plot.

195 **Figure 2.** Maximum likelihood phylogenetic tree based on a nucleotide sequence between the
196 positions 2,872 - 3,152 (S2) of the aligned SARS-CoV-2 genomes.

197 **Figure 3.** Maximum likelihood phylogenetic tree based on 12 WSPs detected across the SARS-
198 CoV-2 genomes.

199

200 **Methods**

201 A total of 1,137 full-length genomes of SARS-CoV-2 sampled from December 2019 to
202 March 25, 2020 (at 2:30 pm) were obtained from Genbank³⁵ and GISAID³⁶ (Supplementary Table
203 S1). Only genomes with high sequencing coverage, intact ORFs (no frameshifts, except that of
204 *nsp12* cistron) and without any indeterminate nucleotide base (indicated by ‘N’s or ambiguous
205 codes) totalizing 767 high quality full-length sequences were effectively analyzed in this study.

206 The genomic data set was aligned using MAFFT-FFT-NS-2³⁷. The calculation of the
207 average number of nucleotide differences per site (nucleotide diversity, π) was conducted in
208 DnaSP v.6³⁸ using a sliding window and step size of 300 and 20 nucleotides, respectively. Sites
209 with gaps alignment were not considered in analysis. The detection of polymorphic sites was
210 conducted using PAUP* v. 4.0³⁹ and MEGA X⁴⁰. Those sites responsible for the segregation of
211 the isolates into two clusters in the ML-trees were referred to as “widely shared polymorphisms”
212 (WSPs), while the remaining nt sites in the virus genomes were designed as “non widely shared
213 polymorphisms” (nWSPs). The WSPs were conventionally indicated by *cistron/gene* name-[nt
214 position within the cistron or gene]. We opted by indicating the nt position within the cistron or
215 gene due to their highly conserved sizes (no gap was introduced during the construction of
216 sequence alignments), in contrast to the full-genomes whose 5’- and 3’-untranslated regions
217 (UTRs) were highly variable in terms of length.

218 Maximum likelihood (ML) phylogenetic trees were constructed using RAxML⁴¹ under the
219 nucleotide substitution model General Time-Reversible with gamma distribution (GTRGAMMA).
220 The branch support for ML-trees based on WSPs-containing and nWSPs segments was assessed
221 with 1,000 and 5,000 bootstrap replications, respectively. All phylogenetic trees were edited using
222 iTOL⁴².

223 References

- 224
- 225 1. Wu, F. *et al.* A new coronavirus associated with human respiratory disease in China. *Nature*
226 **579**, 265–269 (2020).
- 227 2. Coronaviridae Study Group of the International Committee on Taxonomy of Viruses. The
228 species Severe acute respiratory syndrome-related coronavirus: classifying 2019-nCoV and
229 naming it SARS-CoV-2. *Nat. Microbiol.* **5**, (2020).
- 230 3. WHO. WHO Director-General’s remarks at the media briefing on 2019-nCoV on 11
231 February 2020. *WHO website* <https://www.who.int/dg/speeches/detail/who-directo> (2020).
- 232 4. Sawicki, S. G. & Sawicki, D. L. Coronavirus Transcription: A Perspective. in *Current topics*
233 *in microbiology and immunology* vol. 287 31–55 (2005).
- 234 5. de Wilde, A. H., Snijder, E. J., Kikkert, M. & van Hemert, M. J. Host Factors in Coronavirus
235 Replication. in *Assessment & Evaluation in Higher Education* vol. 37 1–42 (Springer
236 International Publishing, 2017).
- 237 6. Kim, D. *et al.* *The architecture of SARS-CoV-2 transcriptome.* (2020) doi:10.1088/1751-
238 8113/44/8/085201.
- 239 7. Peck, K. M. & Luring, A. S. Complexities of Viral Mutation Rates. *J. Virol.* **92**, (2018).
- 240 8. Simmonds, P., Aiewsakun, P. & Katzourakis, A. Prisoners of war — host adaptation and its
241 constraints on virus evolution. *Nat. Rev. Microbiol.* **17**, 321–328 (2019).
- 242 9. Domingo, E., Sheldon, J. & Perales, C. Viral Quasispecies Evolution. *Microbiol. Mol. Biol.*
243 *Rev.* **76**, 159–216 (2012).
- 244 10. Domingo, E. & Perales, C. Viral quasispecies. *PLoS Genet.* **15**, 1–20 (2019).
- 245 11. Zhao, Z. *et al.* Moderate mutation rate in the SARS coronavirus genome and its
246 implications. *BMC Evol. Biol.* **4**, 1–9 (2004).
- 247 12. Gorbalenya, A. E., Enjuanes, L., Ziebuhr, J. & Snijder, E. J. Nidovirales: Evolving the
248 largest RNA virus genome. *Virus Res.* **117**, 17–37 (2006).
- 249 13. Ma, Y. *et al.* Structural basis and functional analysis of the SARS coronavirus nsp14–nsp10
250 complex. *Proc. Natl. Acad. Sci.* **112**, 9436–9441 (2015).
- 251 14. Knipe, D. M. & Howley, P. M. *Fields virology. Viruses and the Lung: Infections and Non-*
252 *Infectious Viral-Linked Lung Disorders* (2013).
- 253 15. Forni, D., Cagliani, R., Clerici, M. & Sironi, M. Molecular Evolution of Human
254 Coronavirus Genomes. *Trends Microbiol.* **25**, 35–48 (2017).
- 255 16. Angelini, M. M., Akhlaghpour, M., Neuman, B. W. & Buchmeier, M. J. Severe Acute
256 Respiratory Syndrome Coronavirus Nonstructural Proteins 3, 4, and 6 Induce Double-
257 Membrane Vesicles. *MBio* **4**, 1–10 (2013).
- 258 17. Cottam, E. M., Whelband, M. C. & Wileman, T. Coronavirus NSP6 restricts
259 autophagosome expansion. *Autophagy* **10**, 1426–1441 (2014).
- 260 18. Gao, Y. *et al.* Structure of RNA-dependent RNA polymerase from 2019-nCoV, a major
261 antiviral drug target. *bioRxiv* 2020.03.16.993386 (2020) doi:10.1101/2020.03.16.993386.
- 262 19. Hoffmann, M. *et al.* SARS-CoV-2 Cell Entry Depends on ACE2 and TMPRSS2 and Is
263 Blocked by a Clinically Proven Protease Inhibitor. *Cell* 1–10 (2020)
264 doi:10.1016/j.cell.2020.02.052.
- 265 20. Li, W. *et al.* Angiotensin-converting enzyme 2 is a functional receptor for the SARS
266 coronavirus. *Nature* vol. 426 450–454 (2003).
- 267
- 268

- 269 21. Matsuyama, S. *et al.* Efficient Activation of the Severe Acute Respiratory Syndrome
270 Coronavirus Spike Protein by the Transmembrane Protease TMPRSS2. *J. Virol.* **84**, 12658–
271 12664 (2010).
- 272 22. Shulla, A. *et al.* A Transmembrane Serine Protease Is Linked to the Severe Acute
273 Respiratory Syndrome Coronavirus Receptor and Activates Virus Entry. *J. Virol.* **85**, 873–
274 882 (2011).
- 275 23. Rossen, J. W. A. *et al.* The Viral Spike Protein Is Not Involved in the Polarized Sorting of
276 Coronaviruses in Epithelial Cells †. *J. Virol.* **72**, 497–503 (1998).
- 277 24. Yang, Y. *et al.* Two Mutations Were Critical for Bat-to-Human Transmission of Middle
278 East Respiratory Syndrome Coronavirus. *J. Virol.* **89**, 9119–9123 (2015).
- 279 25. Pallesen, J. *et al.* Immunogenicity and structures of a rationally designed prefusion MERS-
280 CoV spike antigen. *Proc. Natl. Acad. Sci. U. S. A.* **114**, E7348–E7357 (2017).
- 281 26. Walls, A. C. *et al.* Unexpected Receptor Functional Mimicry Elucidates Activation of
282 Coronavirus Fusion. *Cell* **176**, 1026–1039 (2019).
- 283 27. Walls, A. C. *et al.* Structure, Function, and Antigenicity of the SARS-CoV-2 Spike
284 Glycoprotein. *Cell* **180**, 1–12 (2020).
- 285 28. Sung, S.-C., Chao, C.-Y., Jeng, K.-S., Yang, J.-Y. & Lai, M. M. C. The 8ab protein of
286 SARS-CoV is a luminal ER membrane-associated protein and induces the activation of
287 ATF6. *Virology* **387**, 402–413 (2009).
- 288 29. Chen, C. *et al.* Open Reading Frame 8a of the Human Severe Acute Respiratory Syndrome
289 Coronavirus Not Only Promotes Viral Replication but Also Induces Apoptosis. *J. Infect.*
290 *Dis.* **196**, 405–415 (2007).
- 291 30. Cui, J., Li, F. & Shi, Z.-L. Origin and evolution of pathogenic coronaviruses. *Nat. Rev.*
292 *Microbiol.* **17**, 181–192 (2019).
- 293 31. Parker, M. M. & Masters, P. S. Sequence comparison of the N genes of five strains of the
294 coronavirus mouse hepatitis virus suggests a three domain structure for the nucleocapsid
295 protein. *Virology* **179**, 463–468 (1990).
- 296 32. Huang, Q. *et al.* Structure of the N-Terminal RNA-Binding Domain of the SARS CoV
297 Nucleocapsid Protein. *Biochemistry* **43**, 6059–6063 (2004).
- 298 33. Verheije, M. H. *et al.* The Coronavirus Nucleocapsid Protein Is Dynamically Associated
299 with the Replication-Transcription Complexes. *J. Virol.* **84**, 11575–11579 (2010).
- 300 34. Leung, D. T. M. *et al.* Antibody Response of Patients with Severe Acute Respiratory
301 Syndrome (SARS) Targets the Viral Nucleocapsid. *J. Infect. Dis.* **190**, 379–386 (2004).
- 302 35. Sayers, E. W. *et al.* GenBank. *Nucleic Acids Res.* **47**, D94–D99 (2019).
- 303 36. Shu, Y. & McCauley, J. GISAID: Global initiative on sharing all influenza data – from
304 vision to reality. *Eurosurveillance* **22**, 2–4 (2017).
- 305 37. Katoh, K., Misawa, K., Kei-ichi, K. & Miyata, T. MAFFT: a novel method for rapid
306 multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* **30**, 3059–
307 3066 (2002).
- 308 38. Rozas, J. *et al.* DnaSP 6: DNA Sequence Polymorphism Analysis of Large Data Sets. *Mol.*
309 *Biol. Evol.* **34**, 3299–3302 (2017).
- 310 39. Swofford, D. L. PAUP*: Phylogenetic Analysis Using Parsimony (*and Other Methods).
311 (2002) doi:10.1111/j.0014-3820.2002.tb00191.x.
- 312 40. Kumar, S., Stecher, G., Li, M., Knyaz, C. & Tamura, K. MEGA X: Molecular Evolutionary
313 Genetics Analysis across Computing Platforms. *Mol. Biol. Evol.* **35**, 1547–1549 (2018).
- 314

- 315 41. Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
316 phylogenies. *Bioinformatics* **30**, 1312–1313 (2014).
317 42. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v4: recent updates and new
318 developments. *Nucleic Acids Res.* **47**, W256–W259 (2019).

319 **Acknowledgments**

320 This study was financed in part by the Coordenação de Aperfeiçoamento de Pessoal de Nível
321 Superior - Brasil (CAPES) - Finance Code 001. IJM and RCP were recipients of CNPq and CAPES
322 doctoral fellowships, respectively. DIB was the recipient of a FAPEMIG master fellowship.

323

324 **Author Contributions**

325 ATML designed the bioinformatics analyses. IJM, ATML, RCP, GMS, DIB and FTS conducted
326 the analyses. ATML, IJM, RCP and FTS analyzed data and results. IJM, RCP, FTS and ATML
327 wrote the manuscript. All authors contributed to the content and writing of the Supplementary
328 Information.

329

330 **Competing interest declaration**

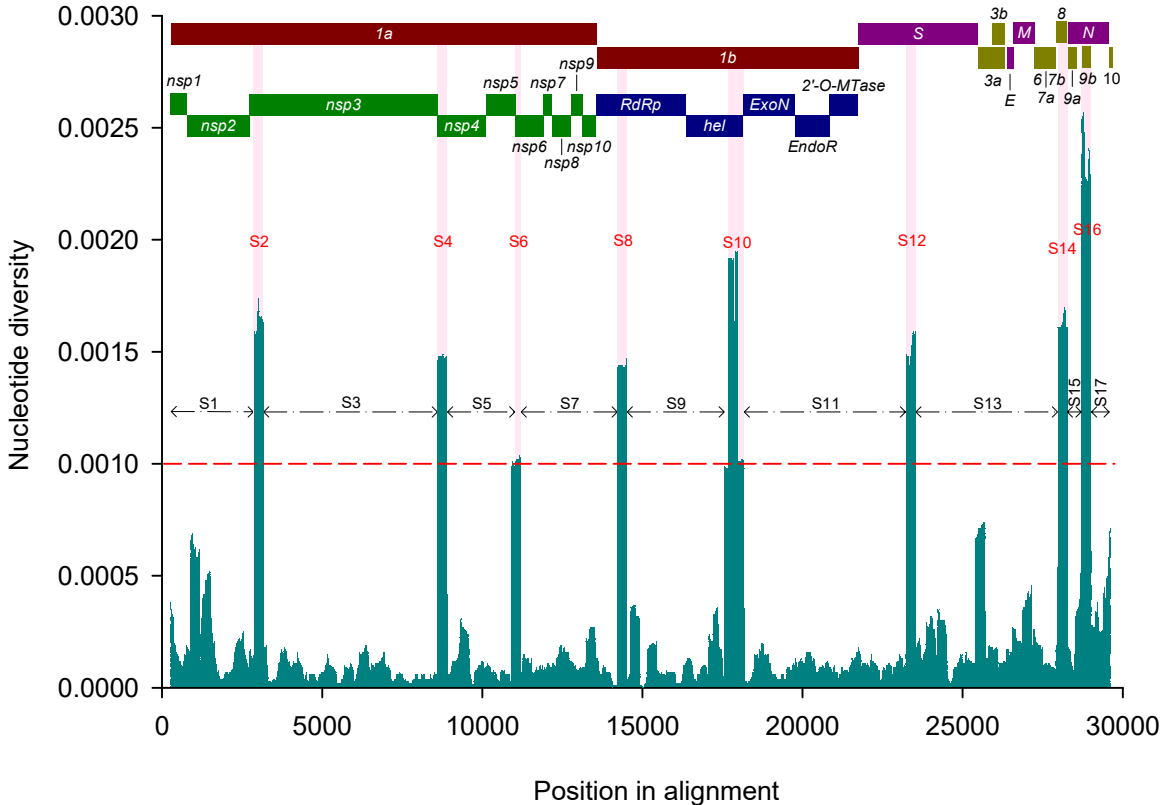
331 The authors declare that they have no competing interests.

332

333 **Additional information**

334 Supplementary information is available for this manuscript.

SARS-CoV-2



Clusters



Major

bioRxiv preprint doi: <https://doi.org/10.1101/2020.04.14.040782>; this version posted April 15, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

Minor

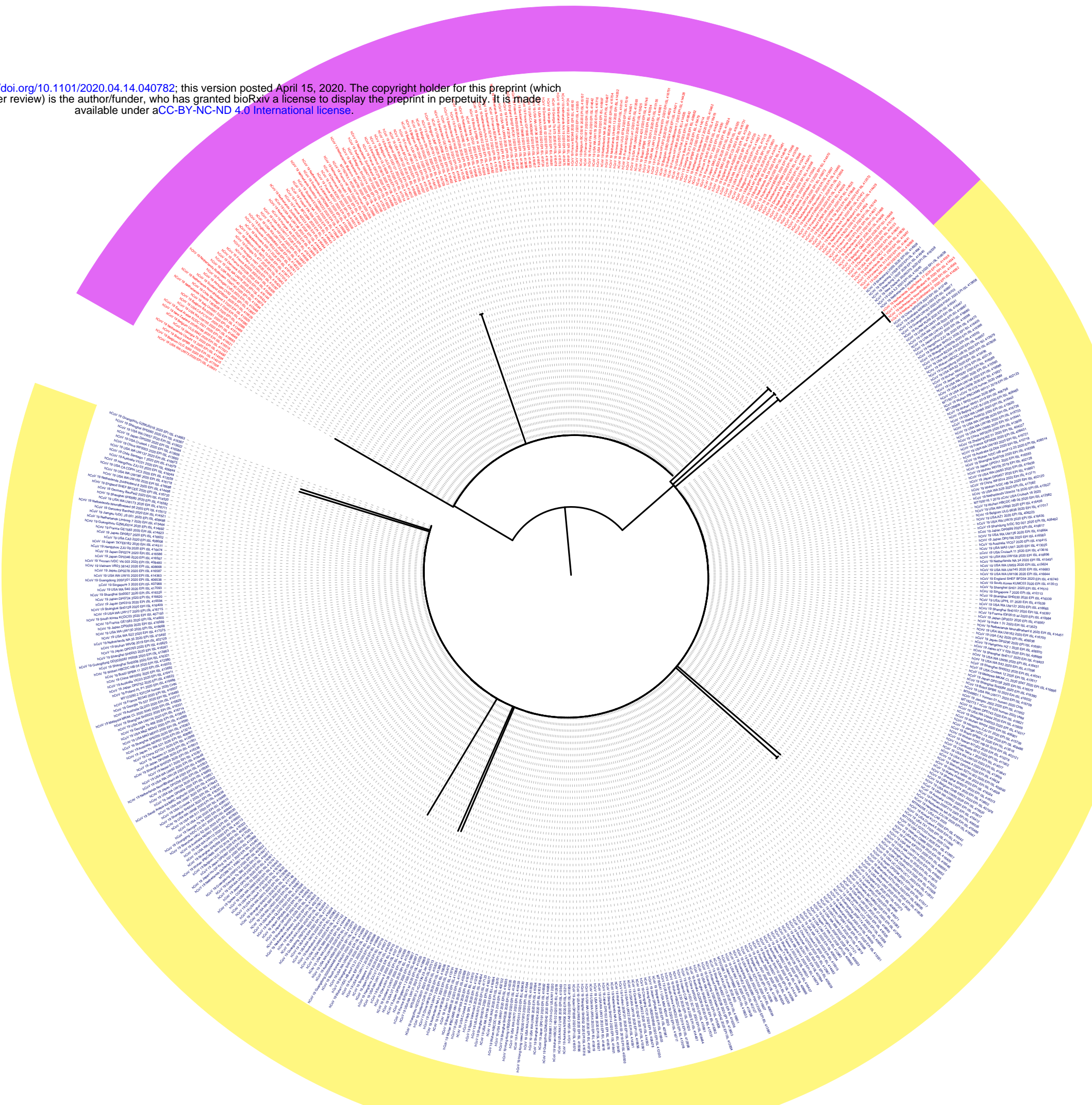
nt mutation



U [N=184]

C [N=409]

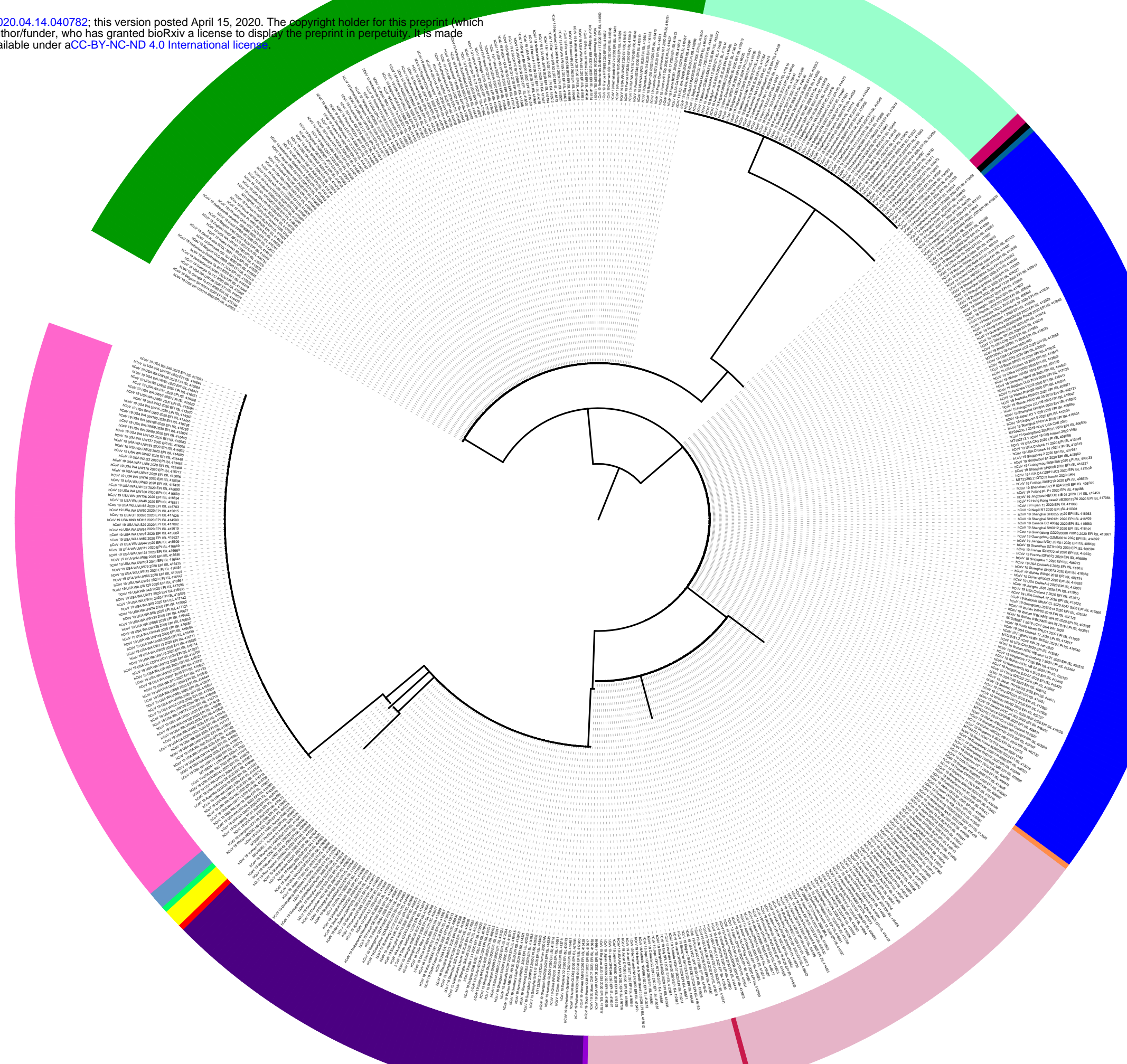
Tree scale: 0.001



Subtypes of SARS-CoV-2

bioRxiv preprint doi: <https://doi.org/10.1101/2020.04.14.040782>; this version posted April 15, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

- I [N=132]
- II [N=122]
- III [N=101]
- IV [N=91]
- V [N=74]
- VI [N=58]
- VII [N=3]
- VIII [N=3]
- IX [N=2]
- X [N=1]
- XI [N=1]
- XII [N=1]
- XIII [N=1]
- XIV [N=1]
- XV [N=1]
- XVI [N=1]



Tree scale: 1