

43 To complement these current classification systems, we consider that haplotypes description
44 and nomenclature could help to better track important mutations that are currently circulating
45 in the world. Identification of SARS-CoV-2 haplotypes aids in understanding the evolution of the
46 virus and may improve our efforts to control the disease.

47 To perform a reasonable analysis of the worldwide temporal and geographical distribution of
48 SARS-CoV-2 haplotypes, we need to take into account the differences in the number of
49 sequenced genomes in months and countries or continents. Thus, we first used a data set of 171
50 461 complete genomes to estimate the worldwide relative frequency of nucleotides in each
51 SARS-CoV-2 genomic position and found nine mutations with respect to the reference genome
52 EPI_ISL_402125 with normalized relative frequencies (NRFp) representing to be present in more
53 than 9 500 000 COVID-19 cases. After that, using a total of 109 953 complete genomes without
54 ambiguous nucleotides from GISAID we performed a phylogenetic analysis and correlated the
55 major branches with SARS-CoV-2 variants which can be classified into five haplotypes or
56 Operational Taxonomic Units (OTUs) based on the distribution of the nine identified nucleotide
57 positions in our NRFp analysis. After that, we analyzed the geographical and temporal worldwide
58 distribution of OTUs normalized by the number of COVID-19 cases. Also, we attempt to correlate
59 these OTUs with patient status, age, and gender information. Finally, we discuss the current
60 hypothesis of the most frequent mutations on protein structure and function. All this
61 information will be continuously updated in our publicly available web-page
62 (<http://sarscov2haplofinder.urp.edu.pe/>).

63 **RESULTS AND DISCUSSION:**

64 **Mutations frequency analysis**

65 The GISAID database contains 171 461 genomes with at least 29 000 sequenced bases; from
66 these, 109 953 genomes do not present ambiguities (as of November 30th). With an alignment
67 of the 171 461 genomes, we performed a normalized relative frequency analysis of each
68 nucleotide in each genomic position (NRFp) (see material and methods for details). This
69 normalization was performed to detect relevant mutations that could appear in regions where
70 few genomes were sequenced (Fig. S1 shows that no correlation exists between the number of
71 cases and the number of sequenced genomes). Using this NRFp analysis, we identified nine
72 positions estimated to be in more than 9 500 000 COVID-19 cases (more than 0.18 NRFp) (Fig.
73 1.A and S2.A) plus many other mutations with NRFp between 0.00-0.18 (Fig. S2.B and S2.C).

74 The nine most frequent mutations (NRFp greater than 0.18) comprise seven non-synonymous
75 mutations, one synonymous mutation, and one mutation in the 5'-UTR region of the SARS-CoV-
76 2 genome (Fig. 1.A). The three consecutive mutations G28881A, G28882A, and G28883C falls at
77 the 5` ends of the forward primer of "China-CDC-N" (Table. S1). Because these three mutations
78 are at the 5` ends, it is unlikely that those mutations greatly affect amplification efficiency. The
79 other six mutations do not fall within regions used by qRT-PCR diagnostic kits (Table. S1). All
80 these nine mutations have been already identified in other studies (Korber et al. 2020, Kern et
81 al. 2020, Pachetti et al. 2020, Yin et al. 2020), although with different frequencies mainly due to
82 the absence of normalization.

83 **OTUs identification**

84 After NRFp analysis, we estimated a maximum likelihood tree using the whole-genome
85 alignment of the 109 953 complete genomes without ambiguities. Then, we associated the
86 branches of the tree with an alignment of the nine positions (241, 1059, 3037, 14408, 23403,

87 25563, 28881, 28882, 28883). We noted that combinations of those nine positions represent
88 five well-defined groups in the tree (Fig. 1.B). Using these combinations, we defined 5
89 haplotypes that allow us to classify more than 97 % of the analyzed genomes (Fig. 1.C), a great
90 part of the remaining not classified genomes are due to the absence of sequencing
91 corresponding to position 241. We named these haplotypes Operational Taxonomic Units
92 (OTUs).

93 OTU_1 was considered the ancestor haplotype due to its identity with the first isolated genomes
94 (EPI_ISL_402125 and EPI_ISL_406801) with characteristic C241, C3037, C14408, and A23403.
95 This OTU_1 comprised genomes with T or C in position 8782 and C or T in 28144. In other
96 analyses, these mutations divide SARS-CoV-2 strains into two lineages. For instance; at the
97 beginning of the pandemic, Tang et al. (2020) showed linkage disequilibrium between those
98 positions and named them as S and L lineages. Rambaut et al. (2020) used these positions to
99 discriminate between their proposed major lineages A and B. Those mutations did not reach the
100 estimated number of 9 500 000 COVID-19 cases, indicating that a small number of these
101 genomes emerged during the pandemic in comparison with other variations.

102 A SARS-CoV-2 isolated on January 25th in Australia is at present the first belonging to OTU_2
103 (Fig. S3). Showing simultaneously four mutations different to OTU_1 (C241T, C3037T, C14408T,
104 and A23403G), OTU_2 is the first group containing the D614G and the P323L mutations in the
105 spike and nsp12 protein, respectively. Korber et al. (2020) analyzed the temporal and geographic
106 distribution of this mutation separating SARS-CoV-2 into two groups, those with D614 and those
107 with G614. Tomaszewski et al. (2020) analyzed the entropy of variation of these two mutations
108 (D614G and P323L) until May. Apparently, OTU_2 is the ancestor of two other OTUs (OTU_3 and
109 OTU_4), as shown in the maximum likelihood tree (Fig. 1.B). OTU_2 is divided into two major
110 branches, one that originates OTU_3 and another more recent branch characteristic from
111 Europe (see below, worldwide geographical distribution of OTUs).

112 On February 16th in the United Kingdom, a SARS-CoV-2 with three adjacent mutations
113 (G28881A, G28882A, and G28883C) (Fig. S3) in N protein was isolated. These three mutations
114 (together with those that characterized OTU_2) define OTU_3. The maximum likelihood tree
115 shows that OTU_4 comes from OTU_2. OTU_4 does not present mutations in N protein; instead,
116 it presents a variation in Orf3a (G25563T). Finally, OTU_5 presents all the mutations of OTU_4
117 plus one nsp2 mutation (C1059T).

118 These nine mutations have been separately described in other reports but, to our knowledge,
119 they have not yet used been used together to classify SARS-CoV-2 haplotypes during the
120 pandemic. The change of relative frequencies of those mutations analyzed individually showed
121 that just in few cases, mutations that define haplotypes described here appear independently
122 (Fig. S4). For example, the four mutations that define OTU_2 (C241T, C3037T, C14408T and
123 A23403G) rarely had been described separately and similarly with mutations that characterize
124 OTU_3 (G28881A, G28882A, G28883C) (Fig. S4). Thus, in this case analysis of haplotypes will be
125 identical results that if we analyzed those mutations independently.

126 The fact that we were able to classify more than 97 % of the complete genomes data set (Fig.
127 1.C) shows that, at least to the present date, this classification system covers almost all the
128 currently known genomic information around the world. Also, most of the unclassified tips
129 appear within a clade allowing us to easily establish their phylogenetic relationships to a
130 haplotype. Thus, at the moment this system can be of practical use to analyze the geographical
131 and temporal distribution of haplotypes during these eleven months of 2020. For convenience

132 we presented table S2 that contains the relation between our identified OTUs and their
133 relationships with pangolin lineages (Rambaut et al. 2020) and GISAID clades (Shu et al. 2017).

134 **Worldwide geographic distribution of OTUs**

135 Using our OTUs classification, we analyzed the worldwide geographic distribution during eleven
136 months of 2020. We began by plotting continental information in the ML tree of the
137 unambiguous complete genomes (Fig. 2.A) and observed some interesting patterns. For
138 instance, all continents contain all OTUs; also, is relatively clear that most isolates belonging to
139 OTU_5 come from North America (Fig. 2.A). Furthermore, the biggest branch of OTU_2 is almost
140 exclusively filled by genomes from Europe, is interesting to note that this branch also contains
141 genomes isolated in the last months analyzed showing its relatively recent appearance (see
142 below, the worldwide temporal distribution of OTUs). However, this approach does not allow us
143 to evaluate continents with less sequenced genomes (Fig. S5.A), such as South America, Oceania,
144 and Africa. Also, it is possible that fine differences can be found in the frequency of one OTU
145 concerning another in each continent. These differences are not observed at this level of
146 analysis.

147 To better analyze which were the most prevalent OTUs in each continent, we analyzed all the
148 complete genomes in the GISAID database (171 461 genomes). In this analysis, we compared
149 the mean of the frequency of OTUs normalized by cases in each continent of six randomly
150 selected groups of genomes (see material and methods for more details).

151 This approach more clearly illustrates that OTU_5 was the most prevalent in North America,
152 followed by OTU_2 and OTU_3, the least prevalent were OTU_1 and OTU_4 (Fig. 2.B). The first
153 genomes in North America belonged to OTU_1 (Fig. S6). Since March, North America was
154 dominated by OTU_5 (Fig. S6). OTU_5 has six of the nine high-frequency genomic variations
155 described (all except those in N protein) (Fig. 1.A).

156 South America presents a greater OTU_3 frequency (Fig. 2.C) that was established in April (Fig.
157 S5). This observation correlates well with studies focused in South America that detect the
158 establishment of D614G mutation at the end of March (mutation presents in OTU_2, OTU_3,
159 OTU_4 and OTU_5) and a high frequency of pangolin lineage B1.1 in Chile and in general in South
160 America that contains the same characteristics mutations that our OTU_3 (Castillo et al. 2020,
161 Franco-Muñoz et al. 2020). Unfortunately, few genomes are reported in South America for
162 September, October, and November (24 genomes in total in the three months), hindering a
163 correct analysis of frequencies in these months. Similarly, OTU_3 was most prevalent in Asia,
164 Oceania, and Africa (Fig. 2.E, 2.F, and 2.G). With other OTUs with least than 0.3 NRFp (Fig. 2.E,
165 2.F, and 2.G). Wu et al. 2020 reports high incidence of mutations that define OTU_3 in
166 Bangladesh, Oman, Russia, Australia and Latvia. At the haplotype level, OTU_3 presents
167 mutations in the N protein that apparently increases the fitness of this group in comparison with
168 OTU_2 (OTU_2 does not present mutations in N) (Fig. 1.A). Thus, four of the six continents
169 analyzed presents an estimation of more than 50 % COVID-19 cases with a SARS-CoV-2 with the
170 three mutations in the N protein. We, therefore, believe that is important to more deeply study
171 if exists positive fitness implications for these mutations.

172 Europe presents an interesting pattern, it follows a similar pattern to South America, Asia,
173 Oceania, and Africa until July (Fig. S6), with OTU_3 as the predominant. Then, in August, OTU_2
174 increased its frequency, and since September OTU_2 is the most prevalent in Europe. This could
175 be caused by the appearance of mutations in the background of OTU_2 (such as those described

176 in Justo et al. 2020) with greater fitness than those of OTU_3 or due to other effects (i.e., founder
177 effects) after the relaxation of lockdown policies.

178 **Worldwide temporal distribution of OTUs**

179 A rooted tree was estimated with the 109 953 genomes data set and labeled by date (Fig. 3.A).
180 Here, we can observe that OTU_1 is mostly labeled with colors that correspond to the first
181 months of the pandemic, expected due to its relation with the first genomes isolated. Clades,
182 where OTU_2, OTU_3, OTU_4, and OTU_5 are the most prevalent, have similar distributions,
183 with representatives mostly isolated since April. The biggest branch of OTU_2 presents a very
184 specific temporal distribution with almost all the genomes isolated from September to
185 November.

186 To gain more insight into these patterns, we estimated the most prevalent OTUs in the world
187 during each month of the pandemic following similar steps that those done for continents (see
188 material and methods for details). In this analysis, we did not consider December and January
189 that present all genomes except one belonging to OTU_1 and mainly from Asia (Fig.S6 and S7).

190 Analysis using the data of February from North America, Europe, and Asia showed that OTU_1
191 continued as the most prevalent in the world but with first isolations of OTU_2, OTU_3, OTU_4,
192 and OTU_5 (Fig. 3.B). Analysis by continents showed that during this month Asia and North
193 America still had higher proportions of OTU_1, but in Europe, a more homogeneous distribution
194 of OTU_1, OTU_2 and OTU_3 was observed (Fig. S6).

195 In March, when the epicenter of the pandemic moved to Europe and North America, but cases
196 were still appearing in Asia, OTU_2, OTU_3, and OTU_5 increased their prevalence but OTU_1
197 remained slightly as the most prevalent during this month (Fig. 3.C). Interestingly, OTU_4
198 remained in relatively low frequencies (Fig. 3.C). This month contains the more homogenous
199 OTUs distribution in a worldwide context, but with some OTUs more prevalent in each continent
200 (Fig. S6).

201 During April, OTU_1 continued its downward while OTU_3 and OTU_5 increased their presence
202 (Fig. 3.D) probably due to its higher representation (compared to March) in several continents
203 such as South America, North America, and Europe (Fig. S6). During this month, Africa showed
204 a high prevalence of OTU_2 (Fig. S6). We also witnessed the establishment of OTU_3 in South
205 America and OTU_5 in North America (Fig. S6).

206 May, June, and July showed a similar pattern, with OTU_3 as the most prevalent due to its high
207 frequencies in South America, Oceania, and Europe (Fig. 3.E, 3.F, 3.G, and S6). North America
208 maintains OTU_5 as the most prevalent and Oceania showed a relatively homogenous pattern.
209 During these months, OTU_2 had intermediate frequencies in all continents resulting in
210 intermediate frequencies all over the world (Fig. 3.E, 3.F, 3.G, and S6). OTU_1 and OTU_4
211 representatives were reported during these months but with very low frequencies.

212 In August and September, we detected a slightly higher frequency of OTU_4 compared to the
213 previous months (Fig. 3H and 3I) with no significant differences with OTU_5. In September in
214 Europe, OTU_3 stopped being the most frequent. Instead, OTU_2 was the most frequent in this
215 month in Europe (Fig. S6). In October and November, OTU_2 has increased its frequency rapidly
216 (Fig. 3.J and 3.H) mainly due to a large number of cases and reported genomes belonging to this
217 OTU_2 in Europe in October and November. Due to the few genomes currently available in

218 GISAID for all continents, except for Europe and North America during November, just these two
219 continents were analyzed in the last month.

220 Also, it is important to mention that, there are not many enough genomes reported for
221 September, October, and November for South America, so during these months OTUs
222 frequencies of this continent were not considered.

223 **Age, Gender and Patient Status relation with OTUs**

224 Relating the distribution of haplotypes according to patient information can help to determine
225 the preference of some OTUs for some characteristics of the patients. Thus, we analyze OTUs
226 distribution according to age, gender, and patient status information available as metadata in
227 the GISAID database.

228 Unfortunately, just 26.11 % of the 171 461 genomes analyzed have age and gender information
229 (Fig. S8). In the case of patient status information, we noted that GISAID categories are not well
230 organized and we had to reclassify the information into three categories; Asymptomatic, Mild,
231 and Severe (Fig. S9.A). Using this classification scheme, we found that 99.14 % (169 979
232 genomes) were not informative, 0.1 % (175 genomes) falls in the Asymptomatic category, 0.33
233 % (562 genomes) in the Mild category and 0.43 % (745 genomes) could be classified as Severe
234 (Fig. S9.B).

235 Using this limited data, we attempt to determine whether any OTU causes an asymptomatic,
236 mild, or severe infection more frequently. We look for significant differences between the
237 relative frequencies of the OTUs in total samples and samples with known patient information.
238 If we found differences, it would mean that some OTU could be more or less related to one type
239 of infection. Here, we analyzed just the month-continent combination with at least 45 genomes
240 with information of one type of infection and at least two times of genomes with any
241 information (for example Asia – February has 58 Asymptomatic genomes and 613 total
242 genomes). Ten combinations meet these criteria, one in the asymptomatic category, one in the
243 mild, and eight in the severe. None of the OTUs frequencies in samples with patient status
244 information were significant different from the frequencies in the total population of the
245 month-continent analyzed (Fig. 4). Thus, we concluded that none of the OTUs are related to an
246 asymptomatic, mild, or severe COVID-19, at least in the populations analyzed.

247 Age information was also analyzed in the same manner. In general, although some differences
248 were detected as significant, those were not consistently maintained between different
249 populations analyzed (Fig. S10.A-J). Furthermore, none difference reaches a p-value less than
250 0.01 (Except for OTU_4 in North America). Since heterogeneity between countries information
251 is possible, we think that these small differences are more likely due to these heterogeneities
252 and we cannot strongly conclude that some age groups are more related to a specific OTU.
253 Additionally, a strong positive correlation between total relative frequencies of OTUs and
254 relative frequencies by age groups in month-continent was found, meaning that those two
255 frequencies are similar in most of the analyzed populations (Fig. S10.K)

256 A similar approach was done using gender information, but in this case, due to the greater
257 quantity of information, we used more restrictive filter parameters. Thus, we selected country-
258 month combinations with at least 250 genomes with male or female information and two times
259 total genomes information (for instance USA – March has 2079 genomes from female patients
260 and 9287 genomes with or without gender information). Again, we did not find OTU's preference
261 for a specific gender (Fig. S11).

262 **Description of the most frequent mutations**

263 C241T

264 The C241T mutation is present in the 5`-UTR region. In coronaviruses, the 5`-UTR region is
265 important for viral transcription (Madhugiri et al. 2014) and packaging (Masters. 2019).
266 Computational analysis showed that this mutation could create a TAR DNA-binding protein 43
267 (TDP43) binding site (Mukherjee and Goswami. 2020), TDP43 is a well-characterized RNA-
268 binding protein that recognizes UG-rich nucleic acids (Kuo et al. 2014) described to regulate
269 splicing of pre-mRNA, mRNA stability and turnover, mRNA trafficking and can also function as a
270 transcriptional repressor and protect mRNAs under conditions of stress (Lee et al. 2011).
271 Experimental studies are necessary to confirm different binding constants of TDP43 for the two
272 variants of 5`-UTR and its *in vivo* effects.

273 C1059T

274 Mutation C1059T lies on Nsp2. Nsp2 does not have a clearly defined function in SARS-CoV-2
275 since the deletion of Nsp2 from SARS-CoV has little effect on viral titers and so maybe
276 dispensable for viral replication (Graham et al. 2005). However, Nsp2 from SARS-CoV can
277 interact with prohibitin 1 and 2 (PBH1 and PBH2) (Cornillez-Ty et al. 2009), two proteins involved
278 in several cellular functions including cell cycle progression (Wang et al. 1999), cell migration
279 (Rajalingam et al. 2005), cellular differentiation (Sun et al. 2004), apoptosis (Fusaro et al. 2003),
280 and mitochondrial biogenesis (Merkwirth and Langer. 2008).

281 C3037T

282 Mutation C3037T is a synonymous mutation in Nsp3; therefore, it is more difficult to associate
283 this change with an evolutionary advantage for the virus. This mutation occurred in the third
284 position of a codon. One possibility is that this changes the frequency of codon usage in humans
285 increasing expression or any other of the related effects caused by synonymous codon change
286 (some of them reviewed in Mauro and Chapel. 2014).

287 C3037T causes a codon change from TTC to TTT. TTT is more frequently present in the genome
288 of SARS-CoV-2 and other related coronaviruses compared to TTC (Gu et al. 2014) but in humans,
289 the codon usage of TTT and TTC are similar (Mauro and Chapel. 2014). The reason why TTT is
290 more frequent in SARS-CoV-2 is unknown but seems to be a selection related to SARS-CoV-2 and
291 not to the host. Another option is genetic drift.

292 C14408T

293 The C14408T mutation changes P323 to leucine in Nsp12, the RNA-dependent RNA polymerase
294 of SARS-CoV2 (Fig. S12.A and B). P323 together with P322 ends helix 10 and generate a turn that
295 is followed by a beta-sheet (Fig. S12.C). Leucine at position 323 could form hydrophobic
296 interactions with the methyl group of L324 and the aromatic ring of F396 creating a more stable
297 variant of Nsp12 (Fig. S12.E). In concordance with this, protein dynamics simulations showed a
298 stability increase of the Nsp12 P323L variant (Chand and Azad. 2020). In the absence of P322,
299 the mutation P323L would probably be disfavored due to the flexibilization of the turn at the
300 end of helix 10. Experimental evidence is necessary to confirm these hypotheses and to evaluate
301 their impact on protein function.

302 A23403G

303 An interesting protein to track is spike protein (Fig. S13.A) due to its importance in SARS-CoV-2
304 infectivity. It has been suggested that the D614G change in the S1 domain that results from the
305 A23403G mutation generates a more infectious virus, less spike shedding, greater incorporation
306 in pseudovirions (Zhang et al. 2020), and higher viral load (Korber et al. 2020).

307 How these effects occur at the structural level remains unclear, although some hypotheses have
308 been put forward: 1) We think that there is no evidence for hydrogen-bond between D614 and
309 T859 mentioned by Korber et al. 2020, distances between D614 and T859 are too long for a
310 hydrogen bond (Fig S13.B), 2) distances between Q613 and T859 (Fig. S13.C) could be reduced
311 by increased flexibility due to D614G substitution, forming a stabilizing hydrogen bond, 3)
312 currently available structures do not show salt-bridges between D614 and R646 as proposed by
313 Zhang et al. 2020 (Fig. S13.D).

314 G25563T

315 Orf3a (Fig. S14.A) is required for efficient in vitro and in vivo replication in SARS-CoV (Castaño-
316 Rodríguez et al. 2018). It has been implicated in inflammasome activation (Siu et al. 2019),
317 apoptosis (Chan et al. 2009), necrotic cell death (Yue et al. 2018) and has been observed in Golgi
318 membranes (Padhan et al. 2007) where pH is slightly acidic (Griffiths and Simons. 1986). Kern et
319 al. 2020 showed that Orf3a preferentially transports Ca²⁺ or K⁺ ions through a pore (Fig S14.B).
320 Some constrictions were described in this pore, one of them formed by the side chain of Q57
321 (Fig. S14.C).

322 Mutation G25563T produces the Q57H variant of Orf3a (Fig. S14.C). It did not show significant
323 differences in expression, stability, conductance, selectivity, or gating behavior (Kern et al.
324 2020). We modeled Q57H mutation and we did not observe differences in the radius of
325 constriction (Fig. S14.C) formed by residue 57 but we observed slight differences in the
326 electrostatic surface due to the ionizability of the histidine side chain (Fig. S14.D).

327 G28881A, G28882A, G28883C

328 N protein is formed by two domains and three disordered regions. The central disordered region
329 named LKR was shown to interact directly with RNA (Chang et al. 2009) and other proteins (Luo
330 et al. 2005), probably through positive side chains; also, this region contains phosphorylation
331 sites able to modulate the oligomerization of N protein (Chang et al. 2013).

332 Mutation G28883C that changes a glycine for arginine at position 204 contributes one more
333 positive charge to each N protein. Mutations G28881A and G28882A produce a change from
334 arginine to lysine. These two positive amino acids probably have a low impact on the overall
335 electrostatic distribution of N protein. However, change from R to K could alter the probability
336 of phosphorylation in S202 or T205. Using the program NetPhosK (Blom et al. 2004), we
337 observed different phosphorylation potential in S202 and T205 between G28881-G28882-
338 G28883 (RG) and A28881-A28882-C28883 (KR) (Fig. S15). Other authors proposed that these
339 mutations could change the molecular flexibility of N protein (Rahman et al. 2020).

340 **CONCLUDING REMARKS:**

341 Here, we present a complete geographical and temporal worldwide distribution of SARS-CoV-2
342 haplotypes from December 2019 to November 2020. We identified nine high-frequency
343 mutations. These important variations (asserted mainly by their frequencies) need to be tracked
344 during the pandemic.

345 Our haplotypes description showed to be phylogenetically consistent, allowing us to easily
346 monitor the spatial and temporal changes of these mutations in a worldwide context. This was
347 only possible due to the unprecedented worldwide efforts in the genome sequencing of SARS-
348 CoV-2 and the public databases that rapidly share the information.

349 Our geographical and temporal analysis showed that OTU_3 is currently the more frequent
350 haplotype circulating in four of six continents (Africa, Asia, Oceania, and South America), result
351 that is in accordance with other studies (Mercatelli et al. 2020) that showed GISAID clade GR
352 (that corresponds to our OTU_3) as the most prevalent in the world; however, they did not
353 report the currently predominance of OTU_2 in Europe (clade G for GISAID). Intriguingly, OTU_3
354 never reached frequencies higher than OTU_5 in North America. In Europe, currently and
355 different from the tendency from May to July, OTU_2 is now much more commonly isolated
356 than OTU_3. Why mutations R203K and G204R have such frequencies in most of the continents,
357 why in North America those mutations were not so successful and why currently Europe is
358 dominated by OTU_2 are open questions. Some studies showed that at the moment there are
359 not mutations that significant increase the fitness of the SARS-CoV-2 (Rasmussen et al. 2020,
360 van Dorp et al. 2020).

361 Although OTU_1 was the only and the most abundant haplotype at the beginning of the
362 pandemic, now its isolation is rare. This result shows an expected adaptation process of SARS-
363 CoV-2. This enunciate does not mean that SARS-CoV-2 is now more infectious or more
364 transmissible.

365 In the next months, these haplotypes description will need to be updated, identification of new
366 haplotypes could be performed by combining the identification of new frequent mutations and
367 phylogenetic inference. We will continue monitoring the emergence of mutations that exceed
368 our proposed cut-off of 0.18 NRFp and this information will be rapidly shared with the scientific
369 community through our web page (<http://sarscov2haplofinder.urp.edu.pe/>). This will also be
370 accompanied by a continuous update of haplotypes information. During the peer-review
371 process o this manuscript, we identify several other mutations near to the cut-off proposed that
372 were reported in Justo et al. 2020.

373 Using information of specific populations we showed no preference for patient's features (age,
374 gender, or type of infection) by OTUs. Thus, mutations that define those haplotypes do not have
375 a relevant impact on the severity of the disease neither are implied preferentially in infections
376 to males, females, or age.

377 Finally, although more studies need to be performed to increase our knowledge of the biology
378 of SARS-CoV-2, we were able to make hypotheses about the possible effects of the most
379 frequent mutations identified. This will help in the development of new studies that will impact
380 vaccine development, diagnostic test creation, among others.

381 **MATERIAL AND METHODS:**

382 **Normalized frequency analysis of each base or gap by genomic position:**

383 To perform the mutation frequency analysis, we first downloaded a total of 171 461 complete
384 and high coverage genomes from the GISAID database (as of November 30th, 2020). This set of
385 genomes was aligned using ViralMSA using default parameter settings, and EPI_ISL_402125
386 SARS-CoV-2 genome from nt 203 to nt 29674 as the reference sequence (Moshiri. 2020, Li.
387 2018). Subalignments corresponding to genomes divided by continent-month combinations was

388 extracted and relative frequencies of each base or gap in each genomic position were calculated
389 ($RF_{p,m-c}$) using a python script. These relative frequencies were multiplied by the number of
390 cases reported in the respective continent-month combination (CN_{m-c}) obtaining an
391 estimation of the number of cases that present a virus with a specific base or gap in a specific
392 genomic position ($RF_p CN_{m-c}$). Finally, we added the $RF_p CN_{m-c}$ of each subalignment and
393 divided it by the total number of cases in the world ($\sum_{m-c} RF_o CN_{m-c1} / TCN_w$). This procedure
394 allows us to obtain a relative frequency normalized by cases of each base or gap in each genomic
395 position (NRF_p). The number of cases of each country was obtained from the European Centre
396 for Disease Prevention and Control: [https://www.ecdc.europa.eu/en/publications-](https://www.ecdc.europa.eu/en/publications-data/download-todays-data-geographic-distribution-covid-19-cases-worldwide)
397 [data/download-todays-data-geographic-distribution-covid-19-cases-worldwide](https://www.ecdc.europa.eu/en/publications-data/download-todays-data-geographic-distribution-covid-19-cases-worldwide). We used the
398 number of cases of countries with at least one genome sequenced and deposited in GISAID
399 database. Also, we just consider in the analysis month-continent combinations with at least 90
400 genomes sequenced.

401 **Phylogenetic tree construction:**

402 Using an alignment of the 109 953 complete, high coverage genomes without ambiguities, we
403 estimated a maximum likelihood tree with Fasttree v2.1.10 with the next parameters: -nt -gtr -
404 gamma -sprlength 1000 -spr 10 -refresh 0.8 -topm 1.5 close 0.75 (Price et al. 2009, Price et al.
405 2010), after the generation of the tree we improved topology using -boot 1000 and the first
406 output tree as an input using -intree option. To generate the rooted tree (against
407 EPI_ISL_402125) we used the R package treeio, and to generate tree figures with continent or
408 date information by tip we used the ggtree package in R (Yu. 2020, Yu et al. 2017).

409 **OTUs determination:**

410 Mutations respect to EPI_ISL_402125 with NRF_p greater than 0.18 were extracted from the
411 alignment of the non-ambiguous data set of 109 953 genomes and were associated with the
412 whole-genome rooted tree using the MSA function from the ggtree package (Yu. 2020, Yu et al.
413 2017) in R. Then, we visually examined to identify the major haplotypes based in these positions,
414 designated as OTUs (Operational Taxonomic Units). Haplotypes identification based in our NRF_p
415 calculation reduced the bias of the different number of genomes sequenced in each continent
416 and each month by integrating the less biased information of the number of cases. Although,
417 other biases are more difficult, if possible, to reduce or eliminate.

418 **Analysis of OTUs geographical distribution:**

419 In this analysis, we randomly separate the genomes into 6 samples of 28 576 genomes each.
420 Genomes in each sample was divided by continents and by months. In these divisions, OTUs
421 relative frequencies were calculated for each OTU in each month-continent combination
422 ($O_n F_{m-c}$). Then, we multiplied these ($O_n F_{m-c}$) frequencies by the number of cases
423 corresponding to the respective month-continent (CN_{m-c}) to obtain an estimation of the
424 number of cases caused by a specific OTU in a respective month-continent ($O_n CN_{m-c}$). After,
425 these products were grouped by continents, and those from the same continent were added
426 and then divided by the total number of cases in the continent analyzed ($\sum_{m-c1} O_n CN_{m-c1} /$
427 TCN_{c1}). Thus, obtaining a frequency normalized by cases for each OTU in each continent. Finally,
428 following this procedure in each sample, we statistically compared the mean of those six
429 samples using the package "ggpubr" in R with the non-parametric Kruskal-Wallis test, and
430 pairwise statistical differences were calculated using non-parametric Wilcoxon test from the
431 same R package. The number of cases of each country was obtained from the European Centre

432 for Disease Prevention and Control: [https://www.ecdc.europa.eu/en/publications-](https://www.ecdc.europa.eu/en/publications-data/download-todays-data-geographic-distribution-covid-19-cases-worldwide)
433 [data/download-todays-data-geographic-distribution-covid-19-cases-worldwide](https://www.ecdc.europa.eu/en/publications-data/download-todays-data-geographic-distribution-covid-19-cases-worldwide). We used the
434 number of cases of countries with at least one genome sequenced and deposited in GISAID
435 database. Also, we just consider in the analysis month-continent combinations with at least 90
436 genomes sequenced.

437 **Analysis of OTUs temporal distribution:**

438 Following a similar procedure used in the geographical analysis, we now grouped the products
439 $O_n CN_{m-c}$ by months, added them, and then divided by the total number of cases in the analyzed
440 month ($\sum_{m1-c} O_n CN_{m1-c}$)/ TCN_{m1} . As in the geographical analysis, the mean of the six
441 samples was statistically compared using the same procedures and with exactly the same
442 considerations of month-continent combinations.

443 **Analysis of age, gender, and patient status with OTUs distribution:**

444 We determine if OTUs have a preference for age or gender, or cause a COVID-19 with a specific
445 severity. For patient status and age information we selected populations with at least 45
446 genomes in the category to analyze and at least two times the total number of genomes (for
447 example Asia – February has 58 asymptomatic genomes and 613 total genomes). For the gender
448 analysis, we selected sample populations with at least 250 genomes in the category to analyze
449 and at least two times the total number of genomes (for example, USA – March has 2 079
450 genomes from female patients and 9287 genomes with or without gender information). In each
451 selected sample we used the total data (all genomes corresponding to that continent-month
452 combination) and the data with category information (for example male, female, asymptomatic,
453 severe, 16-30 years, etc.). We randomly divided these two groups of genomes into three
454 samples and calculated OTUs frequencies. The mean of the frequency of each OTUs was
455 compared between the two groups using the non-parametric Wilcoxon or Kruskal-Wallis
456 statistical test. In the case of age information, the relative frequencies of each OTUs of the total
457 genomes and the genomes with category information were correlated using Spearman
458 correlation. All plots were produced in R using “ggpubr” and ggplot2.

459 **DATA AVAILABILITY:**

460 The data that support the findings of this study comes from the GISAID initiative (Shu and
461 McCaluey. 2017) (gisaid.org). Python and R scripts used in this study are available on request
462 from the corresponding author upon reasonable request.

463 **REFERENCES:**

- 464 1. Blom N, Sicheritz-Pontén T, Gupta R, Gammeltoft S, Brunak S. 2004. Prediction
465 of post-translational glycosylation and phosphorylation of proteins from the
466 aminoacid sequence. Vol. 4, 1633-1649. Proteomics.
- 467 2. Castaño-Rodríguez C, Honrubia J, Gutierrez-Alvarez J, DeDiego M, Nieto-Torres
468 J, Jimenez-Guardeño J,, Regla-Nava J, Fernandez-Delgado R, Verdia-Báguena C,
469 Queralt-Martín M, Kochan G, Perlman S, Aguilera V, Sola I, Enjuanes L. 2018. Role
470 of severe acute respiratory syndrome coronavirus viroporins E, 3a, and 8a in
471 replication and pathogenesis. Vol. 9(3), 1-23. American Society for Microbiology.
- 472 3. Castillo AE, Parra B, Tapia P, Lagos J, Arata L, Acevedo A, Andrade W, Leal G,
473 Tambley C, Bustos P, Fasce R, Fernández J. 2020. Geographical Distribution of

- 474 Genetic Variants and Lineages of SARS-CoV-2 in Chile. Vol. 8:562615. Frontiers in
475 public health.
- 476 4. Chan C, Tsoi H, Chan W, Zhai S, Wong C, Yao X, Chan W, Tsui S, Chan H. 2009. The
477 ion channel activity of the SARS-coronavirus 3a protein is linked to its pro-
478 apoptotic function. Vol. 41, 2232-2239. The International Journal of Biochemistry
479 and Cell Biology.
- 480 5. Chand G and Azad G. 2020. Identification of novel mutations in RNA-dependent
481 RNA polymerases of SARS-CoV-2 and their implications. bioRxiv preprint doi:
482 <https://doi.org/10.1101/2020.05.05.079939>.
- 483 6. Chang C, Chen C, Chiang M, Hsu Y, Huang T. 2013. Transient oligomerization of
484 the SARS-CoV N protein – Implication for virus ribonucleoprotein packaging. Vol.
485 8(5), e65045. PlosONE.
- 486 7. Chang C, Hsu Y, Chang Y, Chao F, Wu M, Huang Y, Hu C, Huang T. 2009. Multiple
487 nucleic acid binding sites and intrinsic disorder of severe acute respiratory
488 syndrome coronavirus nucleocapsid protein implications for ribonucleocapsid
489 protein packaging. Vol. 83(5), 2255-2264. Journal of Virology.
- 490 8. Cornillez-Ty C, Liao L, Yates J, Kuhn P, Buchmeier M. 2009. Severe acute
491 respiratory syndrome coronavirus nonstructural protein 2 interacts with a host
492 protein complex involved in mitochondrial biogenesis and intracellular signaling.
493 Vol. 83(19), 10314-10318. Journal of Virology.
- 494 9. Cuccinotta D and Vanelli M. 2020. WHO declares COVID-19 a pandemic. Vol. 91, 157-
495 160. Acta Biomedica.
- 496 10. Franco-Muñoz C, Álvarez-Díaz D, Laiton-Donato K, Wiesner M, Escandón P, Usme-Ciro J,
497 Franco-Sierra N, Flórez-Sánchez A, Gómez-Rangel S, Rodríguez-Calderón L, Barbosa-
498 Ramirez, J, Ospitia-Baez E, Walteros D, Ospina-Martinez M, Mercado-Reyes M. 2020.
499 Substitutions in Spike and Nucleocapsid proteins of SARS-CoV-2 circulating in South
500 America. Vol. 85, 104557. Infection, Genetics and Evolution.
- 501 11. Fusaro G, Dasgupta P, Rastogi S, Joshi B, Chellappan S. 2003. Prohibitin induces
502 the transcriptional activity of p53 and is exported from the nucleus upon
503 apoptotic signaling. Vol. 278(48), 47853-47861. The Journal of Biological
504 Chemistry.
- 505 12. Graham R, Sims A, Brockway S, Baric S, Denison M. 2005. The nsp2 replicase
506 protein of murine hepatitis virus and severe acute respiratory syndrome
507 coronavirus is dispensable for viral replication. Vol. 79(21), 13399-13411. Journal
508 of Virology.
- 509 13. Griffiths G and Simons K. 1986. The trans Golgi network: sorting at the exit site
510 of the golgi complex. Vol. 234, 438-443. Science.
- 511 14. Gu H, Chu D Peiris M, Poon L. 2020. Multivariate Analyses of Codon Usage of
512 SARS-CoV-2 and other betacoronaviruses. bioRxiv preprint doi:
513 <https://doi.org/10.1101/2020.02.15.950568>.
- 514 15. Justo S, Zapata D, Huallpa C, Landa G, Castillo A, Garavito-Salini R, Uceda-Campos G,
515 Pineda R. 2020. Global geographic and temporal analysis of SARS-CoV-2 haplotypes
516 normalized by COVID-19 cases during the pandemic. bioRxiv preprint doi:
517 <https://doi.org/10.1101/2020.07.12.199414>.

- 518 16. Justo S, Zapata D, Huallpa C, Landa G, Castillo A, Garavito-Salini R, Uceda-Campos G,
519 Pineda R. 2020. Analysis of the Dynamics and Distribution of SARS-CoV-2 Mutations and
520 its Possible Structural and Functional Implications. bioRxiv preprint doi:
521 <https://doi.org/10.1101/2020.11.13.381228>.
- 522 17. Kepler L, Hamins-Puertolas M, Rasmussen D. 2020. Decomposing the sources of SARS-
523 CoV-2 fitness variation in the United States. bioRxiv preprint doi:
524 <https://doi.org/10.1101/2020.12.14.422739>.
- 525 18. Kern D, Sorum B, Hoel C, Sridharan S, Remis J, Toso D, Brohawn S. 2020. Cryo-
526 EM structure of the SARS-CoV-2 3a ion channel in lipid nanodiscs. bioRxiv
527 preprint doi: <https://doi.org/10.1101/2020.06.17.156554>.
- 528 19. Korber B, Fischer W, Gnanakaran S, Yoon H, Theiler J, Abfalterer W, Hengartner
529 N, Giorgi E, Bhattacharya T, Foley B, Hastie K, Parker M, Partridge D, Evans C,
530 Freeman T, de Silva T, McDanal C, Perez L, Tang H, Moon-Walker A, Whelan S,
531 LaBranche C, Saphire E, Montefiori D. 2020. Tracking changes in SARS-CoV-2
532 Spike: evidence that D614G increases infectivity of the COVID-19 virus.
533 <https://doi.org/10.1016/j.cell.2020.06.043>. Cell.
- 534 20. Kuo P, Chiang C, Wang Y, Doudeva L, Yuan H. 2014. The crystal structure of TDP-
535 43 RRM1-DNA complex reveals the specific recognition for UG- and TG-rich
536 nucleic acids. Vol. 42(7), 4712-4722. Nucleic Acids Research.
- 537 21. Lee E, Lee V, Trojanowski J. 2011. Gains or losses: molecular mechanisms of
538 TDP43-mediated neurodegeneration. Vol. 13(1), 38-50. Nature Reviews
539 Neuroscience.
- 540 22. Luo H, Chen Q, Chen J, Chen K, Shen X, Jiang H. 2005. The nucleocapsid protein
541 of SARS coronavirus has a high binding affinity to the human cellular
542 heterogeneous nuclear ribonucleoprotein A1. Vol. 579, 2623-2628. FEBS letters.
- 543 23. Madhugiri R, Fricke M, Marz M, Ziebuhr J. 2014. RNA structure analysis of
544 alphacoronavirus terminal genome regions. Vol. 194, 76-89. Virus Research.
- 545 24. Maitra A, Chawla M, Raheja H, Biswas N, Chakraborti S, Kumar Am Ghosh S,
546 Sarkar S, Patra S, Kumar R, Ghosh T, Chatterjee A, Banu H, Majumdar A,
547 Chinnaswamy A, Srinivasan N, Dutta S, Das S. 2020. Mutations in SARS-CoV-2
548 viral RNA identified in Eastern India: Possible implication for the ongoing
549 outbreak in India and impact on viral structure and host susceptibility. Vol. 45,
550 76. 1-18. Journal of Biosciences.
- 551 25. Masters P. 2019. Coronavirus genomic RNA packaging. Vol. 537, 198-207.
552 Virology.
- 553 26. Mauro V and Chapel S. 2014. A critical analysis of codon optimization in human
554 therapeutics. Vol. 20(11), 604-613. Trends in Molecular Medicine.
- 555 27. Mercateli D and Giorgi F. 2020. Geographic and genomic distribution of SARS-
556 CoV-2 mutations. Vol. 11, 1800. Frontiers in Microbiology.
- 557 28. Merkwirth C and Langer T. 2008. Prohibitin function within mitochondria:
558 essential roles for cell proliferation and cristae morphogenesis. Vol. 1793, 27-32.
559 Biochimica et Biophysica Acta.
- 560 29. Moshiri N. 2020. ViralMSA: Massively scalable reference-guided multiple
561 sequence alignment of viral genomes. btaa743.
562 doi:10.1093/bioinformatics/btaa743. Bioinformatics.

- 563 30. Mukherjee M and Goswami S. 2020. Global cataloguing of variations in
564 untranslated regions of viral genome and prediction of key host RNA binding
565 protein-microRNA interactions modulating genome stability in SARS-CoV-2.
566 bioRxiv preprint doi: <https://doi.org/10.1101/2020.06.09.134585>
567 31. Pachetti M, Marini B, Benedetti F, Giudici F, Mauro E, Storici P, Masciovecchio C,
568 Angeletti S, Ciccozzi M, Gallo R, Zella D, Ippodrino R. 2020. Emerging SARS-CoV-
569 2 mutation hot spots include a novel RNA-dependent-RNA polymerase variant.
570 Vol. 18(179), 1-9. Journal of Translational Medicine.
571 32. Padhan K, Tanwar C, Hussain A, Hui P, Lee M, Cheung C, Malik J, Jameel S. 2007.
572 Severe acute respiratory syndrome coronavirus Orf3a protein interacts with
573 caveolin. Vol. 88, 3067-3077. Journal of General Virology.
574 33. Price M, Dehal P, Arkin A. 2009. FastTree: Computing large minimum-evolution
575 trees with profiles instead of a distance matrix. Vol. 26, 1641. Molecular Biology
576 and Evolution.
577 34. Price M, Dehal P, Arkin A. 2010. FastTree 2 – Approximately maximum-likelihood
578 trees for large alignments. Vol. 5(3), e9490. PloSONE.
579 35. Rahman M, Islam M, Alam A, Islam I, Hoque M, Akter S, Rahaman M, Sultana M,
580 Hossain M. 2020. Evolutionary dynamics of SARS-CoV-2 nucleocapsid protein
581 and its consequences. 1–19. Journal of Medical Virology
582 36. Rajalingam K, Wunder C, Brinkmann V, Churin Y, Hekman M, Sievers C, Rapp U,
583 Rudel T. 2005. Prohibitin is required for RAS-induced RAF-MEK-ERK activation
584 and epithelial cell migration. Vol. 7(8), 837-843. Nature Cell Biology.
585 37. Rambaut A, Holmes E, Hill V, O`Toole A, Hill V, McCrone J, Ruis C, du Plessis L,
586 Pybus O. 2020. A dynamic nomenclature proposal for SARS-CoV-2 lineages to
587 assist genomic epidemiology. 2020. Nature Microbiology.
588 <https://doi.org/10.1038/s41564-020-0770-5>
589 38. Saha O, Hossain M, Rahaman M. 2020. Genomic exploration light on multiple
590 origin with potential parsimony-informative sites of the severe acute respiratory
591 syndrome coronavirus 2 in Bangladesh. Vol. 21, 100951. Gene Reports.
592 39. Shu Y and McCauley J. 2017. GISAID: Global initiative on sharing all influenza data
593 – from vision to reality. Vol. 22(13), 1-3. Euro Surveillance.
594 40. Siu K, Yuen K, Castaño-Rodríguez C, Ye Z, Yeung M, Fung S, Yuan S, Chan C, Yuen
595 K, Enjuanes L, Jin D. 2019. Severe acute respiratory syndrome coronavirus ORF3a
596 protein activates the NLRP3 inflammasome by promoting TRAF3-dependent
597 ubiquitination of ASC. Vol. 33. 8865-8877. The FASEB Journal.
598 41. Sun L, Liu L, Yang X, Wu Z. 2004. Akt binds prohibitin 2 and relieves its repression
599 of MyoD and muscle differentiation. Vol. 117(14), 3021-3029. Journal of Cell
600 Science
601 42. Tang X, Wi C, Li X, Song Y, Yao X, Wu X, Duan Y, Zhang H, Wang Y, Qian Z, Cui J,
602 Lu J. 2020. On the origin and continuing evolution of SARS-CoV-2. Vol. 7, 1012-
603 1023. National Science Review.
604 43. Tomaszewski T, DeVries R, Dong M, Bhatia G, Norsworthy M, Zheng X, Caetano-
605 Anollés G. 2020. New pathways of mutational change in SARS-CoV-2 proteomes

- 606 involve regions of intrinsic disorder important of virus replication and release.
607 Vol. 16, 1-18. Evolutionary Bioinformatics
- 608 44. Van Dorp L, Richard D, Tan C, Shaw L, Acman M, Balloux F. 2020. No evidence for
609 increased transmissibility from recurrent mutations in SARS-CoV-2. Vol.
610 11(5986). Nature Communications.
- 611 45. Wang S, Nath N, Adlam M, Chellappan S. 1999. Prohibitin, a potential tumor
612 suppressor, interacts with RB and regulates E2F function. Vol. 18, 3501-3510.
613 Oncogene.
- 614 46. WHO. 2020. <https://www.who.int/emergencies/diseases/novel-coronavirus-2019>
615 Retrieved on 25 August 2020.
- 616 47. Wu S, Tian C, Liu P, Guo D, Zheng W, Huang X, Zhang Y, Liu L. 2020. Effects of SARS-CoV-
617 2 mutations on protein structures and intraviral protein-protein interactions. Journal of
618 medical virology. <https://doi.org/10.1002/jmv.26597>.
- 619 48. Yin C. 2020. Genotyping coronavirus SARS-CoV-2: methods and implication.
620 Genomics. <https://doi.org/10.1016/j.ygeno.2020.04.016>
- 621 49. Yue Y, Nabar N, Shi C, Kamenyeva O, Xiao X, Hwang I, Wang M, Kehrl J. 2018.
622 SARS-Coronavirus open reading frame-3a drives multimodal necrotic cell death.
623 Vol. 9, 1-15. Cell Death and Disease.
- 624 50. Yu G. 2020. Using ggtree to visualize data on tree-like structures. Vol. 69, 1-18-
625 Current Protocols in Bioinformatics.
- 626 51. Yu G, Smith D, Zhu H, Guan Y, Lam T. 2017. GGTREE: an R package for
627 visualization and annotation of phylogenetic trees with their covariates and
628 other associated data. Vol. 8, 28-36. Methods in Ecology and Evolution.
- 629 52. Zhang L, Jackson C, Mou H, Ojha A, Rangarajan E, Iazard T, Farzan M, Choe H. 2020.
630 The D614G mutation in the SARS-CoV-2 spike protein reduces S1 shedding and
631 increases infectivity. bioRxiv preprint doi:
632 <https://doi.org/10.1101/2020.06.12.148726>.
- 633 53. Zhu N, Zhang D, Wang W, Li X, Yang B, Song J, Zhao X, Huang B, Shi W, Lu R, Niu
634 P, Zhan F, Ma X, Wang D, Xu W, Wu G, Gao G, Tan W. 2020. A novel coronavirus
635 from patients with pneumonia in China, 2019. Vol. 382(8), 727-733. The New
636 England Journal of Medicine.

637 **Competing interests:**

638 The authors declare no competing interests

639 **Acknowledgements:**

640 This manuscript has been released as a pre-print at <https://doi.org/10.1101/2020.07.12.199414>,
641 (Justo et al.)

642 We are very grateful to the GISAID Initiative and all its data contributors, i.e. the Authors from
643 the Originating laboratories responsible for obtaining the specimens and the Submitting
644 laboratories where genetic sequence data were generated and shared via the GISAID Initiative,
645 on which this research is based. Complete acknowledgements of the 171 461 genomes used are
646 available in supplementary file (SF1-SF20).

647 We thank Professor Shaker Chuck Farah (Institute of Chemistry – University of Sao Paulo) for
648 English writing corrections and helpful comments. Also, we thank Professors Aline Maria da Silva

649 (Institute of Chemistry – University of Sao Paulo), Joao Renato Rebello Pinho (Albert Einstein
650 Hospital – Sao Paulo) and PhD(c). Deyvid Amgarten (Albert Einstein Hospital – Sao Paulo) for its
651 helpful comments. To the Ricardo Palma University High-Performance Computational Cluster
652 (URPHPC) managers Gustavo Adolfo Abarca Valdiviezo and Roxana Paola Mier Hermoza at the
653 Ricardo Palma Informatic Department (OFICIC) for their contribution in programs and remote
654 use configuration of URPHPC. To Gladys Arevalo Chong for her figure style suggestions. To the
655 Fundação de Amparo à Pesquisa do Estado de São Paulo (FAPESP) graduate scholarship (To S.J.A)
656 2015/13318-4 (to C. S. F.) and Universidad Ricardo Palma (URP) for APC financing.

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683 **Figure 1. Five haplotypes (or OTUs) based in nine positions can classify 97 % of the**
684 **genomes.** A) Table showing haplotype of each OTU, regions, and aminoacids changes caused
685 by these mutations. B) Rooted tree of 109 953 SARS-CoV-2 complete and non-ambiguous
686 genomes associated with an alignment of nine genomic positions (241, 1059, 3037, 14408,
687 23403, 25563, 28881, 28882, 28883) showing a good correlation between haplotypes (OTUs)
based in these nine positions. Tips of the tree where colored based in the OTU. C) Bar diagram
showing OTUs distribution of the genomes (0 correspond to unclassified genomes).

688 **Figure 2. By cases normalized continent distribution of OTUs shows OTU_3 as the most**
689 **prevalent in four of six continents.** A) Unrooted tree of complete non-ambiguous genomes,
tips were colored according to OTUs, and points in each tip were colored according to the
690 continent. B-G) Boxplots of normalized relative frequencies of OTUs in each continent from
691 December 2019 to November 2020 (B, North America; C, South America; D, Europe; E, Asia;
F, Oceania; G, Africa).

692 **Figure 3. By cases normalized temporal distribution of OTUs showed OTU_3 as the most**
693 **prevalent until September.** A) Rooted tree of complete non-ambiguous genomes showing
temporal distribution. Tips were colored by OTUs and points in each tip were colored
694 according to the collection date. B-E) Boxplot of OTUs global distribution in each month (B,
695 February; C, March; D, April; E, May; F, June; G, July; H, August; I, September; J, October; K,
November).

696 **Figure 4. OTUs are not related to the COVID-19 severity.** A-J) Ten different sample
697 populations were analyzed, none of the OTUs frequencies shows significative differences
698 between the total samples and samples taken from genomes with patient status information.
699 Boxplots showed the distribution of three samples, total frequencies are showed in grey and
frequencies from samples with patient status information are colored according the category
(green, asymptomatic; blue, mild; red, severe).

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714 **Supplemental Figures captions:**

715 **Table 1. Region of primers binding and amplification of nine diagnostic tests for SARS-CoV-**
716 **2.**

717 **Table 2. Comparison between different nomenclatures of SARS-CoV-2 lineages.**

718 **Figure S1. Number of genomes sequenced by region is not correlated to the number of**
719 **cases in the same region.** Each point in the plot represents a month-continent combination.
720 There are continents with a high-number of cases but low number of sequenced genomes
721 and inversely, there are continents with relatively few cases but with a large number of
722 sequenced genomes.

723 **Figure S2. Normalized Relative Frequency of each nucleotide by position (NRF_p).** The
724 frequency of each nucleotide in each position was normalized by the number of cases in each
725 continent-month pairs to reduce the bias produced by the different number of sequenced
726 genomes in different months and different continents. In A, Labels are showed for NRF_p
727 greater than 0.18. B and C showed different scales of positions with less than 0.18 NRF_p.

728 **Figure S3. Temporal distribution by day, continent, and OTUs.** Each point in the plot
729 represents one of the 171 461 SARS-CoV-2 genomes analyzed. Points are colored depending
730 on the OTU. Y-axis divides the points in continent and each column represents a day from
731 December 16 to July 23.

732 **Figure S4. Global NRF_p of the nine most frequent mutations by month.** Mutations that
733 define OTU_2 (C241T, C3037T, C14408T, A23403G) showed very similar frequencies
734 indicating that genomes with three, two or one of these mutations are rare. The same for
735 mutations that define OTU_3 (G28881A, G28882A, G28883C). Mutations that define OTU_4
736 (C1059T) and OTU_5 (G25563T) have similar but not identical distributions.

737 **Figure S5. Month and continent distribution of the 171 461 SARS-CoV-2 genomes analyzed.**
738 A) Bars represent genome count in each continent analyzed. Europe and North America are
739 overrepresented in the database. B) Bars in B represent genomes count by month. March,
740 April and October are the best represented months. Bars are labeled by percentage and
741 below by the exact counts.

742 **Figure S6. Temporal distribution by month, continent, and OTU.** Each point in the plot
743 represents a genome and is colored depending on OTU. Points are grouped by continent (Y-
744 axis) and month (x-axis). We saw how haplotypes populations changes during time; for
745 example, OTU_1 seems the most common during the first months (December, January, and
746 February).

747 **Figure S7. Distribution of OTUs in January.** Bar plot of a count of complete genomes isolated
748 in January and deposited in the GISAID database. Most of these genomes belonging to
749 OTU_1, a small fraction corresponds to unclassified genomes and one to OTU_2

750 **Figure S8. Approximately 74 % of the genomes in GISAID database does not have gender**
751 **information.** The plot shows gender distribution of the 171 461 SARS-CoV-2 genomes
752 analyzed. Bars represent genomes count in Male, Female or unknown categories.

753 **Figure S9. More than 90 % of the genomes in the GISAID database does not have an**
754 **informative description of patient status.** A) Table showing which GISAID categories were
755 recategorized in the Asymptomatic, mild or severe categories. All the other genomes were
756 classified as non-informative. B) Distribution of 171 461 genomes in patient status categories
(Asymptomatic, Mild, Severe or No informative).

745 **Figure S10. Age groups are not robustly related to OTUs.** A-J) Ten populations were selected
746 to analyze if OTUs frequencies in an age group is significant different to OTUs frequencies
747 in the total population. None OTU showed a repetitive preference for an age group in the
748 populations analyzed, boxplots are colored by age groups, all means frequencies in the total
749 population (ns, $p>0.05$; *, $0.05>p>0.01$; **, $0.01>p>0.005$; ?, not analyzed). K) Correlation
750 between relative frequencies of OTUs in a specific age group with OTUs frequencies in the
751 whole population. Spearman correlation showed an R value of 0.94 meaning a positive
752 correlation that supports the conclusion that no significant differences exist between OTUs
753 frequencies in age groups compared to the whole population.

751 **Figure S11. OTUs do not have preference for males or females.** A-K) Boxplots of OTUs
752 frequencies from female populations compared to OTUs frequencies in the whole
753 population. None significant difference was observed. L-V) The same as A to K but whole
754 population compared to male populations. Again, no significant differences were observed.
755 Concluding that OTUs do not show gender preferences.

755 **Figure S12. P323L could impact the stability of Nsp12 without disturbing its overall structure.** A) Structure of RNA-dependent RNA polymerase complex (PDB ID: 6YYT). Chains
756 (Nsp12, Nsp7, Nsp8, RNA) are distinguished by colors. Helix 10, Beta-sheet 3, Turn 10-3, and
757 P323 also are differentially colored. B) Structure in A rotated 90 degrees. C) Zoom of the red
758 box in B showed P322 and P323 in the center of Turn 10-3. D) Turn 10-3 with side chains of
759 P323, L324, and F396 in sphere representation to highlight the distance between side chains
760 of P323 and L324. E) P323 in D was computationally replaced by L323. Now, distances
761 between the methyl group of leucine are shorter with L323.

760 **Figure S13. Structural hypotheses about D614G mutation in Spike protein.** A) Structure of
761 the open state of Spike trimer (PDB ID: 6YVB) colored by domains. B) Distances between side
762 chains of two possible rotamers of D614 (1`-D614 and 2`-D614) and T859. Except for 1`-D614
763 and carbonyl group of T859, the other distances seems to be large to form a hydrogen bond.
764 C) Distances between side chains Q613 and T859. These distances are also large to form
765 hydrogen bonds. D) R646 points to the opposite side of D614 showing that there is no salt
766 bridge. B, C, and D show electron density maps of the side chains of the labeled residues.

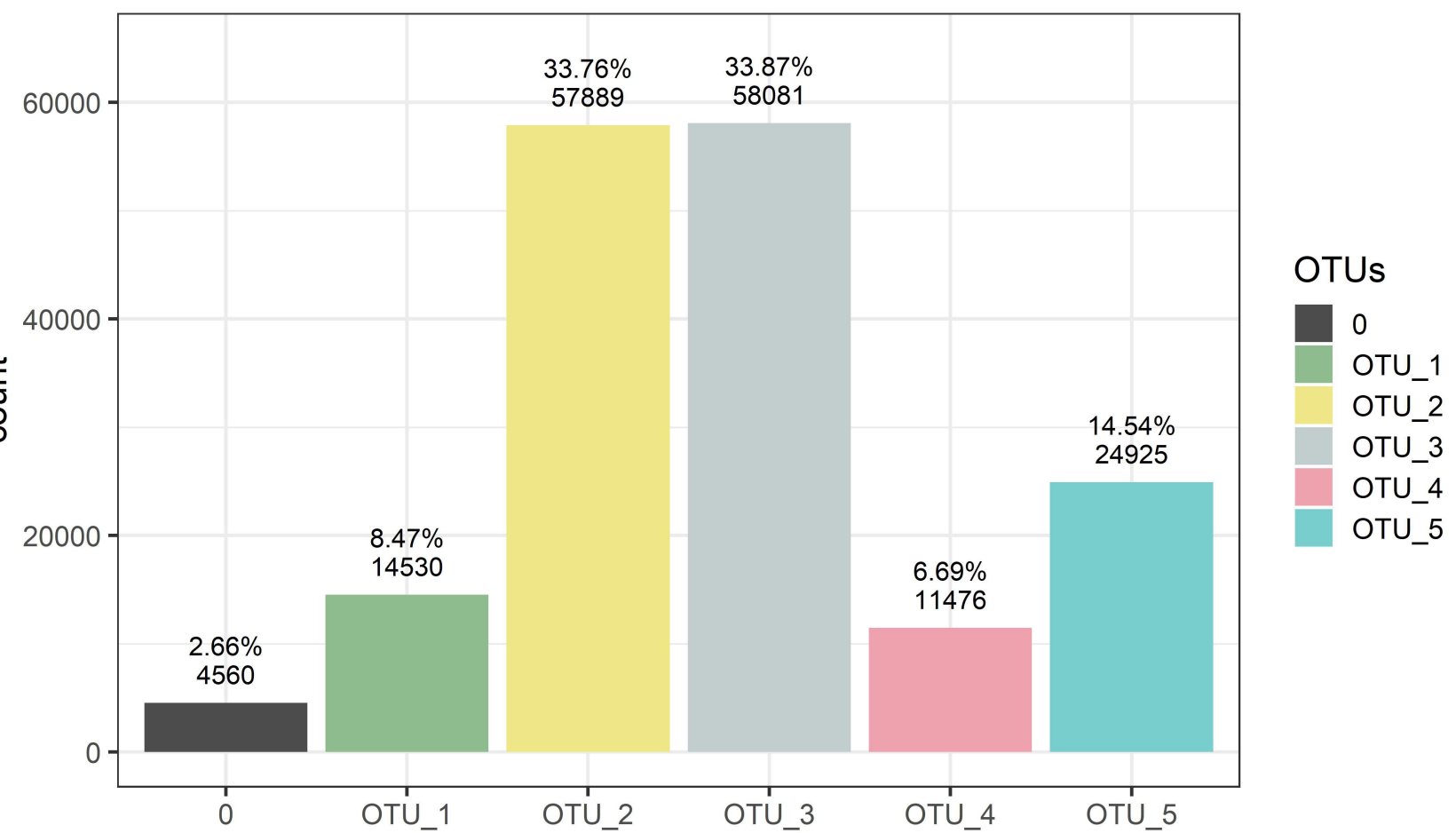
765 **Figure S14. Orf3a Q57H does not modify pore constriction distances but electrostatics distribution.** A) Structure of the Orf3a dimer (PDB ID: 6XDC) colored by domains. The right
766 of A shows the same structure but in an upper view. B) Orf3a showing the central pore, in
767 the red box the section corresponding to the fifth pore constriction. C) zoom of the red box
768 in B, above we showed Q and H variants superposed. Below we show a transversal cut of the
769 pore near to the fifth. The pore radius in two variants is similar. D) Electrostatic surface maps
770 of Q57 and H57 variants in two different pHs (7 and 6). Residues Q57 and H57 are shown in
771 stick representations to point the fifth constriction. We show a slightly more positive region
772 at the height of the fifth constriction.

770 **Figure S15. Mutants in R203 and G204 of Nucleocapsid generate differences in Phosphorylation potential on S202 and T205.** Bar plot showing the phosphorylation
771 potential calculated in NetPhosK for the 4 possible nucleocapsid variants. We can see that
772 phosphorylation potential by PKC is lower for RG than for KR in S202. On the other hand,
773 T205 has greater phosphorylation potential by an unspecific kinase (unsp) in RG than in KR.
774 Phosphorylation in S202 and T205 by unsp or PKC respectively is apparently not affected by
775 these mutations.

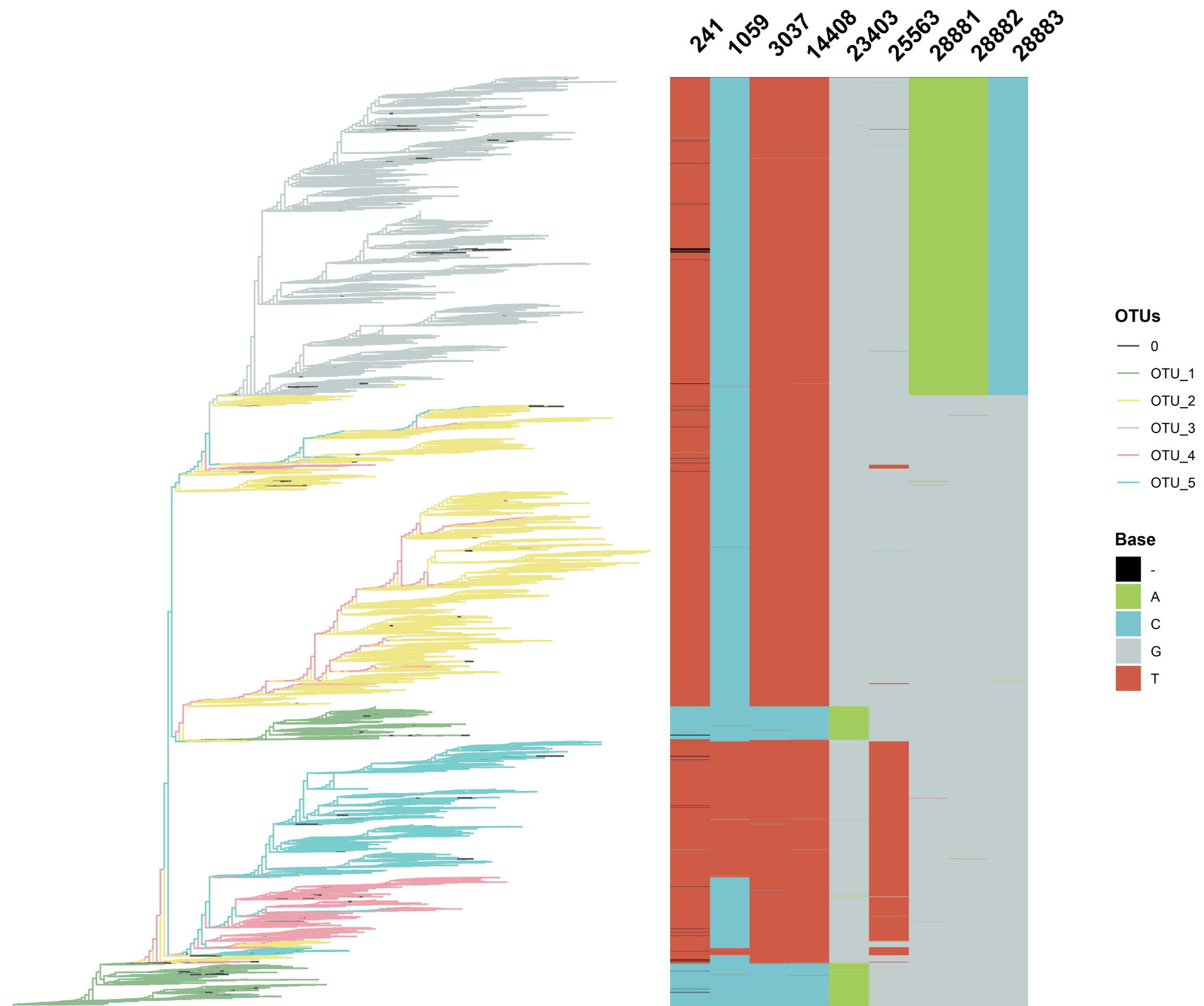
A

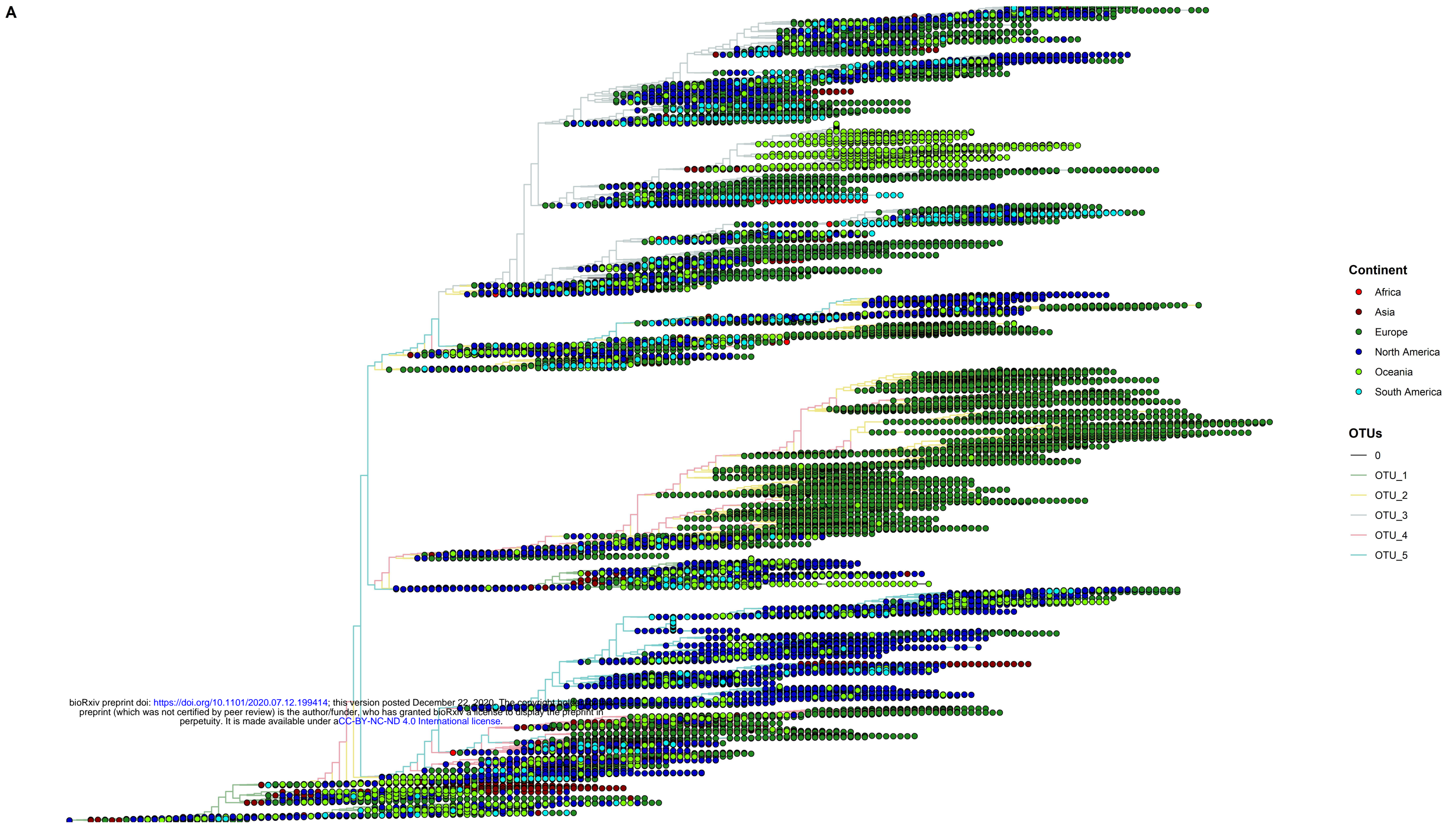
Position	OTU_1	OTU_2	OTU_3	OTU_4	OTU_5	Region	AA_change
241	C	T	T	T	T	5'UTR	''
1059	C	C	C	C	T	Nsp2	T85I
3037	C	T	T	T	T	Nsp3	Syn
14408	C	T	T	T	T	Nsp12	P323L
23403	A	G	G	G	G	S	D614G
25563	G	G	G	T	T	Orf3a	Q57H
28881	G	G	A	G	G	N	R203K
28882	G	G	A	G	G	N	R203K
28883	G	G	C	G	G	N	G204R

C

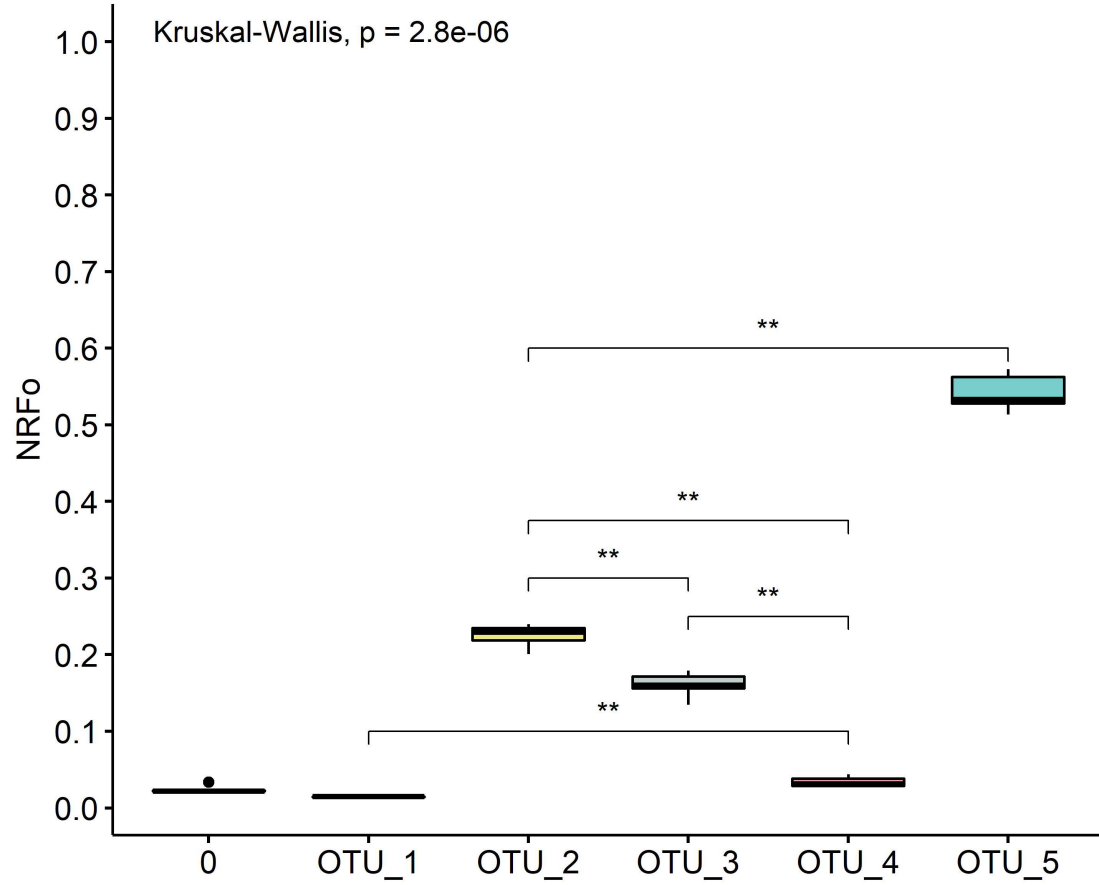


B

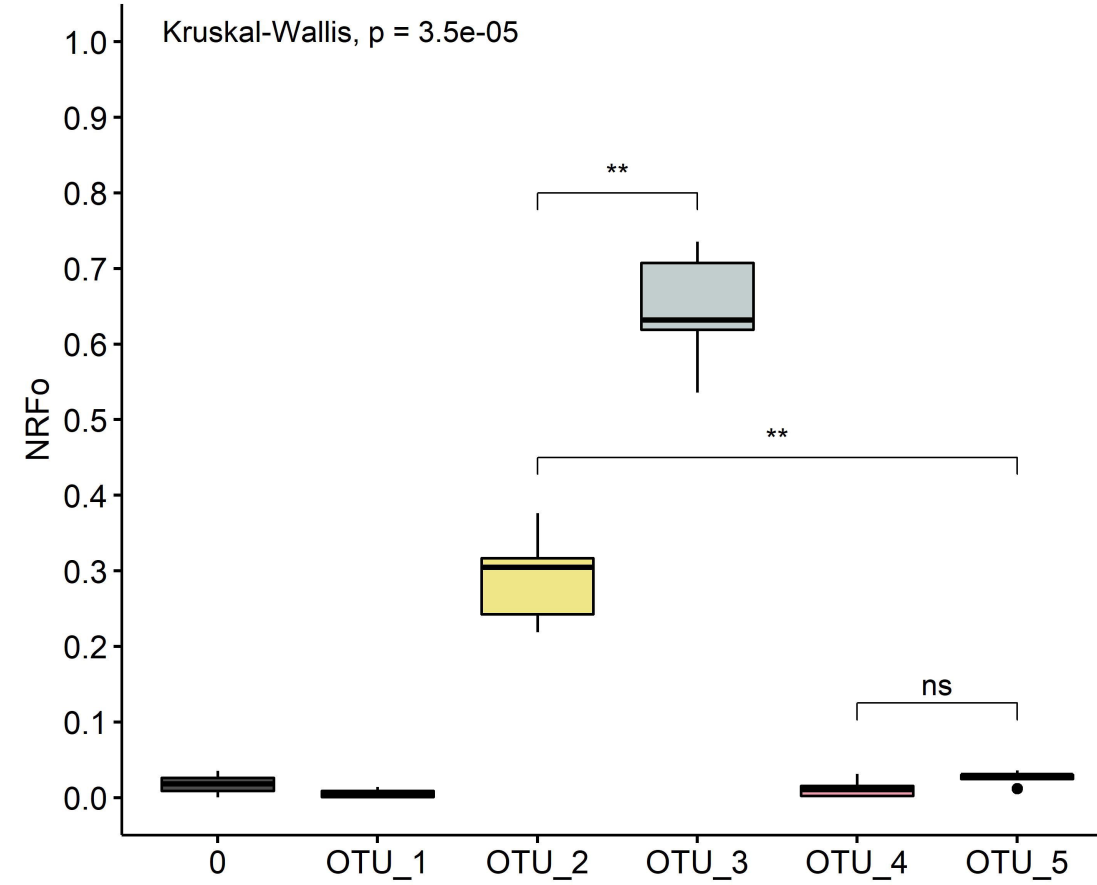




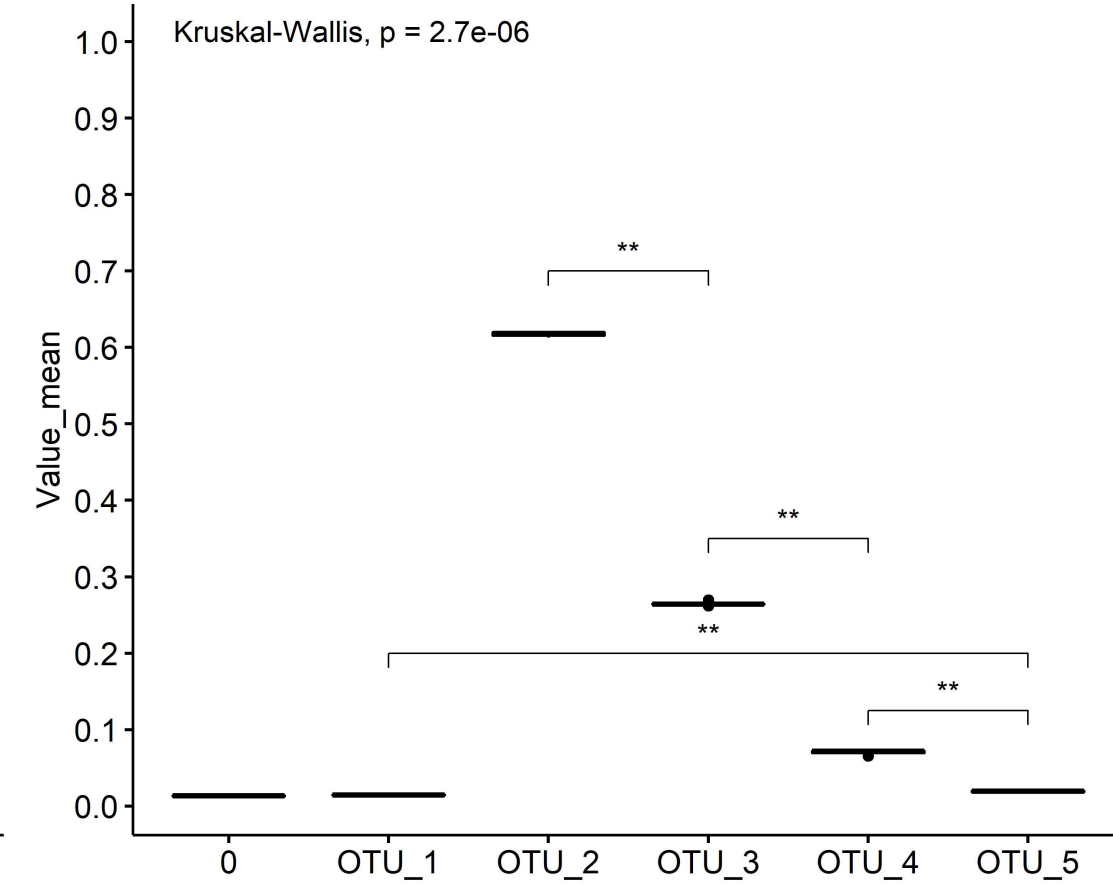
B North America



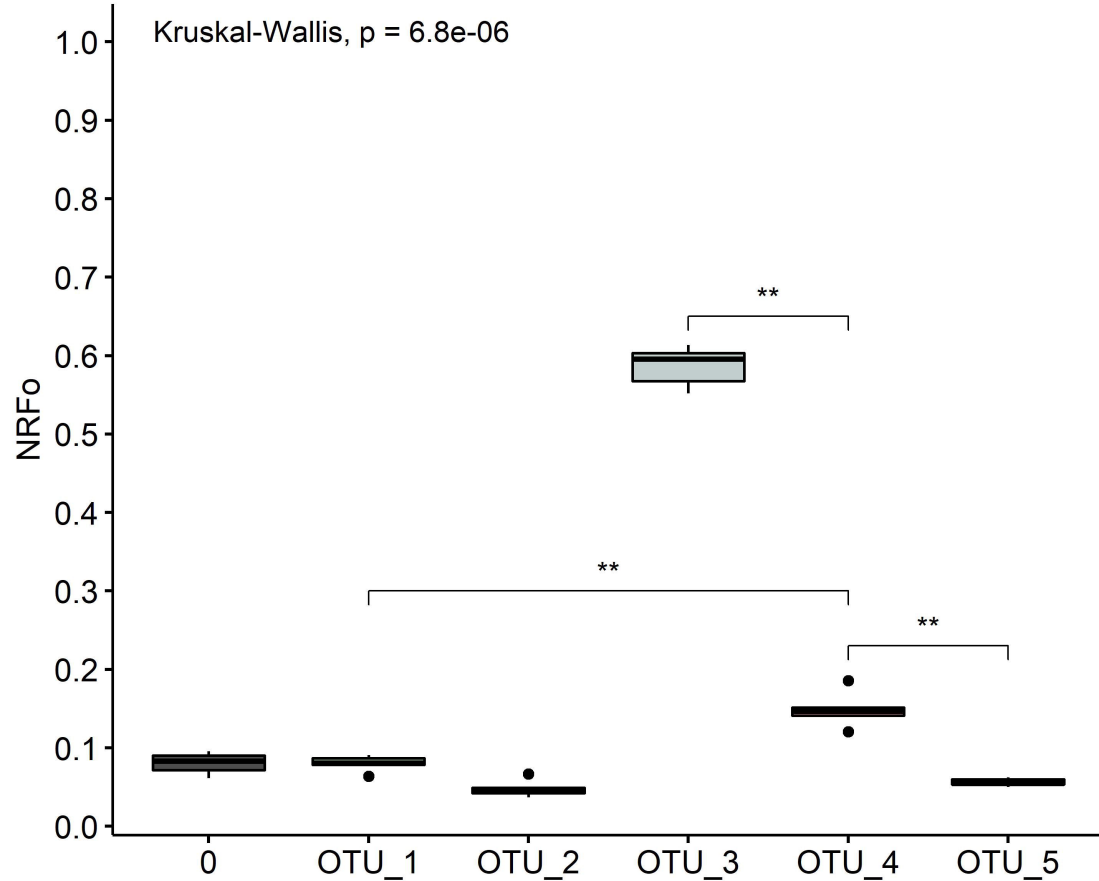
C South America



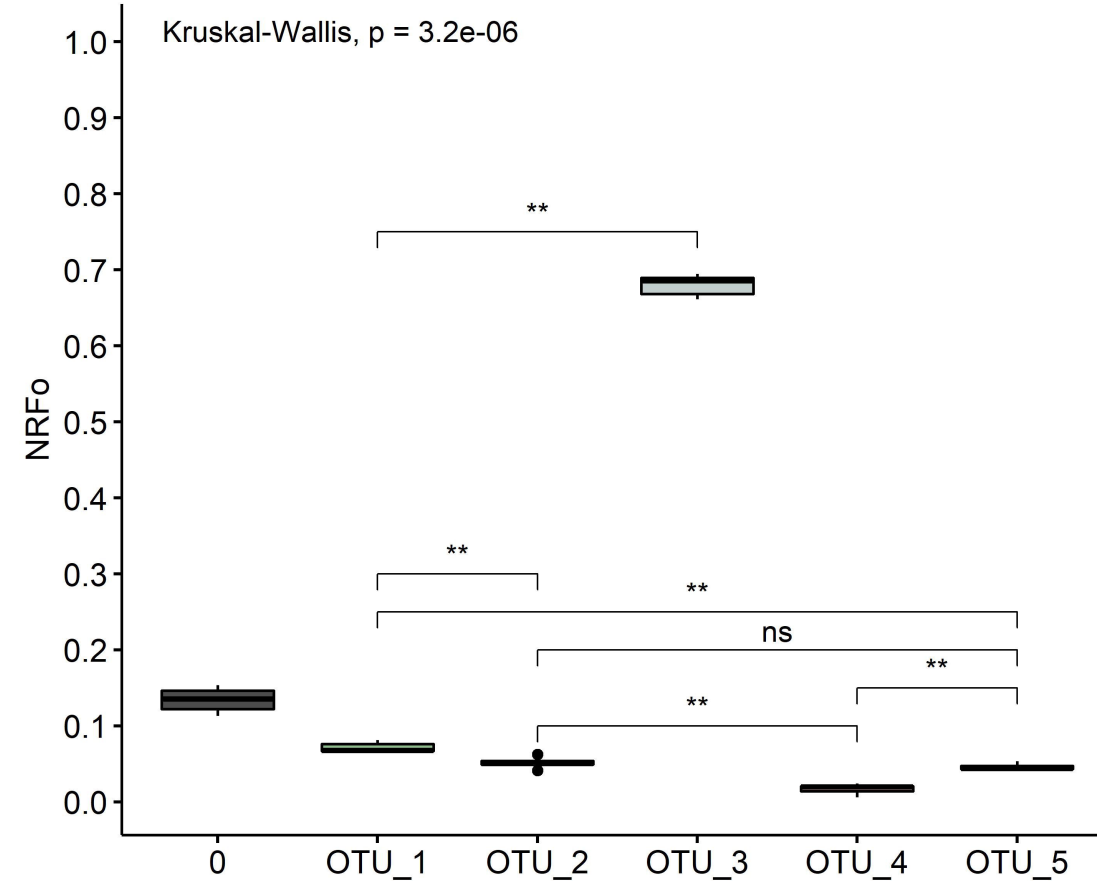
D Europe



E Asia



F Oceania



G Africa

