

1 **Global Geographic and Temporal Analysis of SARS-CoV-2 Haplotypes Normalized by COVID-**
2 **19 Cases during the First Seven Months of the Pandemic**

3 **Santiago Justo Arévalo^{1,2,*}, Daniela Zapata Sifuentes^{1,±}, César Huallpa Robles^{3,±}, Gianfranco Landa**
4 **Bianchi^{1,±}, Adriana Castillo Chávez^{1,±}, Romina Garavito-Salini Casas^{1,±}, Guillermo Uceda-Campos^{2,4}**
5 **Roberto Pineda Chavarría¹.**

6 1.- Universidad Ricardo Palma, Facultad de Ciencias Biológicas, Lima – Perú.

7 2.- Universidade de Sao Paulo, Instituto de Química, Departamento de Bioquímica, São Paulo - Brasil

8 3.- Universidad Nacional Agraria La Molina, Facultad de Ciencias, Lima – Perú.

9 4.- Universidad Nacional Pedro Ruiz Gallo, Facultad de Ciencias Biológicas, Lambayeque - Perú

10 ± These authors contributed equally to this work.

11 * Corresponding author: santiago.justo@urp.edu.pe

12

13 **ABSTRACT:**

14 Since the identification of SARS-CoV-2, a large number of genomes have been sequenced with
15 unprecedented speed around the world. This marks a unique opportunity to analyze virus
16 spreading and evolution in a worldwide context. However, currently, there is not a useful
17 haplotype description to help to track important and globally scattered mutations. Also,
18 differences in the number of sequenced genomes between countries and/or months make it
19 difficult to identify the emergence of haplotypes in regions where few genomes are sequenced
20 but a large number of cases are reported. We proposed an approach based on the normalization
21 by COVID-19 cases of relative frequencies of mutations using all the available data to identify
22 major haplotypes. Thus, we can use a similar normalization approach to tracking the global
23 temporal and geographic haplotypes distribution in the world. Using 48 776 genomes, we
24 identify 5 major haplotypes based on 9 high-frequency mutations. Normalized global geographic
25 and temporal analysis is presented here highlighting the current importance of nucleocapsid
26 mutations (R203K, G204R) above the highly discussed D614G in spike protein. Also, we analyzed
27 age, gender, and patient status distribution by haplotypes, but scarce and not well-organized
28 information about this is publicly available. For that, we create a web-service to continuously
29 update our normalized analysis of mutations and haplotypes, and to allow researchers to
30 voluntarily share patient status information in a well-organized manner to improve analyses and
31 making possible monitor the emergence of mutations and/or haplotypes with patients
32 preferences or different pathogenic features. Finally, we discuss currently structural and
33 functional hypotheses in the most frequently identified mutations.

34 **INTRODUCTION:**

35 COVID-19 was declared a pandemic by the World Health Organization on March 11th 2020¹,
36 with around 23 million cases and 800 thousand of deaths around the world², quickly becoming
37 the most important health concern in the world. Several efforts to produce vaccines, drugs, and
38 diagnostic tests to help in the fight against SARS-CoV-2 are being mounted in a large number of
39 laboratories all around the world.

40 Since the publication on January 24th of the first complete genome sequence of SARS-CoV-2 from
41 China³, thousands of genomes have been sequenced in a great number of countries on all 5
42 continents and were made available in several databases. This marks a milestone in scientific
43 history and gives us an unprecedented opportunity to study how a specific virus evolves in a
44 worldwide context. As of July 30, 2020, the GISAID database⁴ contained 48 776 genomes with at
45 least 29 000 sequenced bases.

46 At the moment, some analysis has been performed to identify SARS-CoV-2 variants around the
47 world, most of them on a particular group of genomes and/or at the beginning of the pandemic
48 using limited datasets. In March 2020 two major lineages were proposed based in position 8782
49 and 28144 using a data set of 103 genomes⁵ which was followed by a particularly interesting
50 proposal that identified the same major lineages (named A and B) and others sublineages⁶.

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

55 To perform a reasonable analysis of the worldwide temporal and geographical distribution of
56 SARS-CoV-2 haplotypes, we need to take into account the differences in the number of
57 sequenced genomes in months and countries-continent. Thus, we first used a data set of 48776
58 complete genomes to estimate the worldwide relative frequency of nucleotides in each SARS-
59 CoV-2 genomic position and found nine positions with normalized relative frequencies (NRF_p)
60 greater than 0.1 and lesser than 0.9. After that, using a total of 19486 complete genomes with
61 any ambiguous nucleotide position from GISAID we performed a phylogenetic analysis and
62 correlated the major branches with SARS-CoV-2 variants which can be classified into five
63 haplotypes or Operational Taxonomic Units (OTUs) based on the distribution of the nine
64 identified nucleotide positions in our NRF_p analysis. After that, we analyzed the geographical
65 and temporal worldwide distribution of OTUs normalized by the number of COVID-19 cases.
66 Also, we attempt to correlate these OTUs with patient status, age, and gender information.
67 Finally, we discuss the current hypothesis of the most frequent mutations on protein structure
68 and function.

69 **RESULTS AND DISCUSSION:**

70 **Mutations frequency analysis**

71 The GISAID database contains around 48 776 genomes with at least 29 000 sequenced bases
72 and from these 19 486 genomes does not contain any ambiguity (as of July 30th). With an
73 alignment of the 48 776 genomes, we performed a normalized relative frequency analysis of
74 each nucleotide in each genomic position (NRF_p) (see material and methods for details), this
75 normalization was performed to reduce the bias due that the number of sequenced genomes in
76 continents and months are not correlated with the number of cases in these continents and
77 months. Using these NRF_p analyses, we identified 9 positions with greater than 0.1 and less than
78 0.9 NRF_p (Fig. 1.A and S1.A) plus many other positions with frequencies between 0.900-0.995
79 and 0.005-0.100 (Fig. S1.B and S1.C).

80 The nine most frequent mutations (NRF_p between 0.1 and 0.9) are comprised of seven non-
81 synonymous mutations, one synonymous mutation and one mutation in the 5'UTR region of the
82 SARS-CoV-2 genome (Fig. 1.A). All these mutations have been already identified in other
83 studies^{7,8,9,10}, although with different frequencies.

84 **OTUs identification**

85 After NRF_p analysis, we estimated a maximum-likelihood tree using the whole-genome
86 alignment of the 19 846 complete genomes without ambiguities. Then, we associated the main
87 branches of the whole-genome tree with an alignment of the 9 positions (241, 1059, 3037,
88 14408, 23403, 25563, 28881, 28882, 28883) and noted that combinations of those 9 positions

89 represent 5 well-defined groups in the tree (Fig. 1.B). Using these combinations, we defined 5
90 haplotypes that allow us to classified 96.5 % of the analyzed genomes (Fig. 1.C), a great part of
91 the remaining not classified genomes are due to the absence of sequencing corresponding to
92 position 241. We named these haplotypes Operational Taxonomic Units (OTUs) and numbered
93 them according to proximity to the root.

94 We were able to clearly track the mutations that originated each of these OTUs. OTU_1 is the
95 ancestor haplotype with characteristic C241, C3037, C14408, and A23403. This OTU_1
96 comprised genomes with T or C in position 8782 and C or T in 28144. In other analyses, these
97 mutations divide SARS-CoV-2 strains into two lineages. For instance; at the beginning of the
98 pandemic, Tang et al (2020) shows linkage disequilibrium between those positions and named
99 them as S and L lineages. Rambaut et al (2020) used these positions to discriminate between
100 their proposed major lineages A and B. After, seven months of pandemic NRF_p of T and C in
101 positions 8782 and 28144, respectively, are not in the range of 0.1-0.9, indicating a small
102 quantity of these genomes presented during the pandemic in comparison with other variations.

103 A SARS-CoV-2 isolated on February 20 was the first belonging to OTU_2 (Fig. S2) that shows
104 simultaneously four mutations different to OTU_1 (C241T, C3037T, C14408T, and A23403G). We
105 can note in the phylogenetic tree that in the transition between the first clade and OTU_2 some
106 unclassified tips were showed. These could be genomes containing some of these four
107 mutations (C241T, C3037T, C14408T, A23403G) but not all, representing intermediate steps in
108 the formation of this haplotype. OTU_2 is the first group containing the D614G mutation in the
109 spike protein. Korber et al. 2020 analyzed the temporal and geographic distribution of this
110 mutation separating SARS-CoV-2 populations into two groups, those with D614 and those with
111 G614.

112 Almost at the same time (February 24), SARS-CoV-2 with three adjacent mutations (G28881A,
113 G28882A, and G28883C) (Fig. S2) in N protein was isolated. These three mutations characterize
114 OTU_3. The maximum likelihood tree shows that OTU_4 comes from OTU_2. OTU_4 does not
115 present mutations in N protein, instead, it presents a variation in Orf3a (G25563T). Finally,
116 OTU_5 presents all the mutations of OTU_5 plus one Nsp2 mutation (C1059T).

117 These 9 mutations have been separately described in other reports but, to our knowledge, they
118 have not yet used been used together to classify SARS-CoV-2 haplotypes during the pandemic.
119 The fact that we were able to classify 96.5 % of the complete genomes data set (Fig. 1.C) shows
120 that, at least to the present date, this classification system covers almost all the currently known
121 genomic information around the world. Also, most of the unclassified tips appear within a clade
122 allowing us to easily establish their phylogenetic relationships to a haplotype. Thus, at the
123 moment this system can be of practical use to analyze the geographical and temporal
124 distribution of haplotypes during these seven months of 2020.

125 It is highly likely that during the next months, some of these OTUs will disappear and others will
126 appear when new mutations in these “parental” OTUs become fixed in the population. Thus,
127 methodologies to actively update circulating haplotypes on a real-time basis need to be
128 proposed. We propose that the best strategy will be to continually monitor the appearance of
129 new haplotypes by tracking mutations that exceed a fixed NRF_p in the world (to allow tracking
130 relevant medical mutations) and associating these mutations to a phylogenomic tree to confirm
131 its phylogenetical relevance, we will perform this task at least one time per month and update
132 this information in our website.

133 **Worldwide geographic distribution of OTUs**

134 Using our OTUs classification, we analyzed the worldwide geographic distribution during the first
135 seven months of 2020. We began by plotting continental information in the unrooted tree of
136 the unambiguous complete genomes (Fig. 2.A) and observed some interesting patterns. For
137 instance, all continents contain all OTUs; however, it is relative clearly that most isolates belonging
138 to OTU_5 come from North America (Fig. 2.A). This approach does not allow us to evaluate
139 continents with less sequenced genomes (Fig. S4), such as South America, Oceania, and Africa;
140 also, it is possible that fine differences can be found in the frequency of one OTU with respect
141 to another in each continent. These differences are not observed at this level of analysis.

142 To better analyze which were the most prevalent OTUs in each continent, we analyzed all the
143 complete genomes in the GISAID database (48 776 genomes). In this analysis, we compared the
144 mean of the frequency of OTUs normalized by cases in each continent of six randomly selected
145 groups of genomes (see material and methods for more details).

146 This approach more clearly illustrates that OTU_5 was the most prevalent in North America,
147 followed by OTUs 2, 3, the least prevalent were OTU_1 and OTU_4 (Fig. 2.B). First genomes in
148 North America belonged to OTU_1 (Fig. S4). March and April were dominated by OTU_5, but in
149 June OTU_3 seems to have similar counts to OTU_5 (Fig. S4). OTU_5 has 6 of the 9 high-
150 frequency genomic variations described (all except those in N protein) (Fig. 1.A).

151 South America presents a greater OTU_3 frequency (Fig. 2.C) that was established in April (Fig.
152 S4). Unfortunately, few genomes were sequenced in South America in May, June, and July (44
153 genomes in total in the three months), hindering a correct analysis of frequencies in these
154 months. Similarly, OTU_3 was most prevalent in Europe, Africa, and Asia (Fig. 2.D, 2.E, and 2.G).
155 Followed by OTU_2 in Europe and Africa, and by OTU_1 in Asia (Fig. 2.D, 2.E, and 2.G). At the
156 haplotype level, OTU_3 present mutations in N protein that apparently increase the fitness of
157 this group in comparison with OTU_3 (OTU_3 does not present mutations in N) (Fig. 1.A). We,
158 therefore, believe that is important to more deeply study the biological implications of these
159 mutations in N protein.

160 Oceania presents a more homogeneous distribution of OTUs, with OTU_1 in slightly higher but
161 statistically significant frequency among other OTUs (Fig. 2.F). The analysis of Oceania is in part
162 biased due to the great percentage of genomes without information of position 241 (in the 5`
163 UTR region), hindering unambiguously classification of several sequenced in Australia.

164 **Worldwide temporal distribution of OTUs**

165 A rooted tree was estimated with the 19 846 genomes data set and labeled by date (Fig. 3.A).
166 Here we can clearly follow the evolution beginning with OTU_1 at the base of the tree (mostly
167 labeled with colors that correspond to the first months of the pandemic). Clades, where OTU_2,
168 4, and 5 are the most prevalent, have intermediate temporal distribution (mostly late February
169 up to late April). OTU_3 has a similar distribution pattern to OTU_2, 4, and 5 but with more
170 representatives isolated in May, June, and July.

171 To gain more insight into these patterns, we estimated the most prevalent OTUs during each
172 month of the pandemic following similar steps that those done for continents (see material and
173 methods for details). In this analysis, we did not consider December and January that present
174 genomes just belonging to OTU_1 mainly from Asia (Fig.S4 and S5).

175 Analysis using the data of February from North America, Europe, and Asia showed that OTU_1
176 continues as the most prevalent in the world but with the presence of OTU_2, 3, 4, and 5 (Fig.

177 3.B). Analysis by continents showed that during this month Asia and North America still had
178 higher proportions of OTU_1, but in Europe, a more homogeneous distribution of OTUs 2-5 was
179 observed (Fig. S4).

180 In March, when the epicenter of pandemic moves to Europe and North America, but cases were
181 still appearing in Asia, OTU_2, 3, and 5 increased its prevalence but OTU_1 remained as the most
182 prevalent during this month (Fig. 3.C). Interestingly OTU_5 remained in relatively low
183 frequencies (Fig. 3.C). Apparently, this month contains the more homogenous OTUs distribution
184 in a worldwide context, but with some OTUs more prevalent in each continent (Fig. S4).

185 During April, OTU_1 continued its downward while OTU_3 and 5 increased its presence (Fig. 3.D)
186 probably due to its higher representation (compared to March) in several continents such as
187 South America, North America, and Europe (Fig. S4). During this month, Africa showed a high
188 prevalence of OTU_2 (Fig. S4). We also witnessed the apparent establishment of OTU_3 in South
189 America and Europe and OTU_5 in North America (Fig. S4).

190 May showed the current tendency of OTU_1 declining and OTU_3 increasing; OTU_2 and 5 were
191 presented in similar frequencies between OTU_1 and 3. OTU_4 maintains its relatively low
192 frequency (Fig. 3.E). From this month, South America reported very few isolated genomes and
193 we cannot consider this continent to the analysis of this and follow months (Fig. S4).

194 The last months analyzed (June and July) presented frequencies distributions very similar to
195 May, showing OTU_3 as the more frequent currently, but with OTU_2 maintaining its frequency,
196 unlike OTU_5 that showed lower frequencies in June when compared with May and July. The
197 current high prevalence of OTU_3 in the world and the observation that also from June in North
198 America, its frequency is rapidly increasing highlights the importance of tracking and study
199 mutations in Nucleocapsid that characterize this OTU.

200 **Age, Gender and Patient Status distribution of OTUs**

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

205 Unfortunately, just 33.65 % of the 48 776 genomes analyzed have age and gender information
206 (Fig. S6) and 4 108 genomes contain some information about the patient status (Fig. S7.B).
207 Distribution of OTUs between age or gender categories did not show any well-defined
208 preference. The distribution of OTUs in different categories was very similar (Fig. 4.A, B, and C).

209 In the case of patient status analysis, we noted that GISAID categories are not well organized
210 and we had to reclassify the information into four categories, Not Informative, Asymptomatic,
211 Mild, and Severe (Fig. S7.A). Using this classification scheme, we found that 55.82 % (2 293
212 genomes) falls in the Not Informative category, 37.22 % (1529 genomes) in the Mild category
213 and just 2.31 % (95 genomes) and 4.65 % (191 genomes) could be classified as Asymptomatic
214 and Severe, respectively (Fig. S7.B).

215 We analyzed the group distribution in the three informative categories (Asymptomatic, Mild and
216 Severe) and found that isolates from patients with mild symptoms presented a relatively
217 homogeneous distribution, with percentages between 27.7 % and 12.1 % from all five OTUs. The
218 severe category was also relatively homogeneous with OTU_1 being the least prevalent (7.9 %).

219 Conversely, 75.8 % (72 of the 95) of the genomes classified as Asymptomatic belong to OTU_1
220 (Fig. 4.D).

221 However, we have to interpret these observations with extreme caution since most of the
222 genomes from asymptomatic patients that belong to OTU_1 was isolated in Asia in February
223 (Fig. S8.A) during a short period of three days (Fig. S8.B). Other genomes in the asymptomatic
224 category belong to other OTUs and were isolated in different months and different continents
225 (Fig. S8.B). Thus, we currently require more robust information to obtain a better-defined
226 distribution of asymptomatic cases, as well as more and better-organized information related
227 to patient status and characteristics to improve our analyses in OTUs distribution related to this
228 data.

229 For this reason, we have created a web page that, in addition to assigning haplotypes to
230 genomes that users can freely upload and make openly available information on the global
231 geographic and temporal distribution of SARS-CoV-2 haplotype in an interactive way, allows
232 researchers from all over the world to contribute voluntarily by offering correctly organized
233 information on the characteristic of the patient (age, gender, condition (symptoms),
234 comorbidities) to improve the analysis and monitor the possible appearance of haplotypes with
235 certain preferences that can help in improving treatments for patients.

236 **Description of the most frequent mutations**

237 C241T

238 The C241T mutation is present in the 5' UTR region. In coronaviruses, the 5'UTR region is
239 important for viral transcription¹¹ and packaging¹². Computational analysis showed that this
240 mutation could create a TAR DNA-binding protein 43 (TDP43) binding site¹³, TDP43 is a well-
241 characterized RNA-binding protein that recognizes UG-rich nucleic acids¹⁴ described to regulate
242 splicing of pre-mRNA, mRNA stability and turnover, mRNA trafficking and can also function as a
243 transcriptional repressor and protect mRNAs under conditions of stress¹⁵. Experimental studies
244 are necessary to confirm different binding constants of TDP43 for the two variants of 5' UTR and
245 its *in vivo* effects.

246 C1059T

247 Mutation C1059T lies on Nsp2. Nsp2 does not have a clearly defined function in SARS-CoV-2
248 since the deletion of Nsp2 from SARS-CoV has little effect on viral titers and so maybe
249 dispensable for viral replication¹⁶. However, Nsp2 from SARS-CoV can interact with prohibitin 1
250 and 2 (PBH1 and PBH2)¹⁷, two proteins involved in several cellular functions including cell cycle
251 progression¹⁸, cell migration¹⁹, cellular differentiation²⁰, apoptosis²¹, and mitochondrial
252 biogenesis²².

253 C3037T

254 Mutation C3037T is a synonymous mutation in Nsp3, therefore, is more difficult to associate this
255 change to an evolutionary advantage for the virus. This mutation occurred in the third position
256 of a codon, one possibility is that this, change the frequency of codon usage in humans
257 increasing expression or any other of the related effects caused by synonymous codon change
258 (some of them reviewed²³).

259 C3037T causes a codon change from TTC to TTT. TTT is more frequently present in the genome
260 of SARS-CoV-2 and other related coronaviruses compared to TTC²⁴ but in humans, the codon

261 usage of TTT and TTC are similar²³. The reason why TTT is more frequent in SARS-CoV-2 is
262 unknown but seems that is a selection related to SARS-CoV-2 and not by the host. Another
263 option is simply genetic drift.

264 C14408T

265 The C14408T mutation changes P323 to leucine in Nsp12, the RNA-dependent RNA polymerase
266 of SARS-CoV2 (Fig. 5.A and B). P323, along with P322 end helix 10, and generate a turn preceding
267 a beta-sheet (Fig. 5.C). Leucine at position 323 could form hydrophobic interactions with the
268 methyl group of L324 and the aromatic ring of F396 creating a more stable variant of Nsp12 (Fig.
269 5.E). Protein dynamics simulations showed an increase in stability of the Nsp12 P323L variant²⁵.
270 In the absence of P322, the mutation P323L would probably be disfavored due to the
271 flexibilization of the turn at the end of helix 10. Experimental evidence is necessary to confirm
272 these hypotheses and to evaluate its impact on protein function.

273 A23403G

274 An interesting protein to track is spike protein (Fig. 6.A) due to its importance in SARS-CoV-2
275 infectivity. It has been suggested that the D614G change in the S1 domain that results from the
276 A23403G mutation generates a more infectious virus, less spike shedding, greater incorporation
277 in pseudovirions²⁶, and higher viral load⁷.

278 How these effects occur at the structural level remains unclear, although some hypotheses have
279 been put forward: 1) We think that there is no evidence for hydrogen-bond between D614 and
280 T859 mentioned by Korber et al. 2020, distances between D614 and T859 are too long for a
281 hydrogen bond (Fig 6.B), 2) distances between Q613 and T859 (Fig. 6.C) could be reduced by
282 increased flexibility due to D614G substitution, forming a stabilizing hydrogen bond, 3) currently
283 available structures do not show salt-bridges between D614 and R646 as proposed by Zhang et
284 al. 2020 (Fig. 6.D).

285 G25563T

286 Orf3a (Fig. 7.A) is required for efficient *in vitro* and *in vivo* replication in SARS-CoV²⁷, has been
287 implicated in inflammasome activation²⁸, apoptosis²⁹, necrotic cell death³⁰ and has been
288 observed in Golgi membranes³¹ where pH is slightly acidic³². Kern et al. 2020 showed that Orf3a
289 preferentially transports Ca⁺² or K⁺ ions through a pore (Fig 7.B) of in which one constriction is
290 formed by the side chain of Q57 (Fig.7.C).

291 Mutation G25563T produces a Q57H variant of Orf3a (Fig. 7.C) that did not show significant
292 differences in expression, stability, conductance, selectivity, or gating behavior⁸. We modeled
293 Q57H mutation and we did not observe differences in the radius of constriction (Fig. 7.C) formed
294 by aminoacid 57 but we observed slight differences in the electrostatic surface due to the
295 ionizability of the histidine side chain (Fig. 7.D).

296 G28881A, G28882A, G28883C

297 N protein is formed by two domains and three disordered regions. The central disordered region
298 named LKR was shown to interact directly with RNA³⁵ and other proteins³⁶, probably through
299 positive side chains; also, this region contains phosphorylation sites able to modulate the
300 oligomerization of N protein³⁷.

301 Mutation G28883C that introduces an arginine at position 204 contributes one more positive
302 charge to each N protein. Mutations G28881A and G28882A produce a change from arginine to

303 lysine, these two positive amino acids probably have a low impact on the overall electrostatic
304 distribution of N protein. However, change from R to K in this position could change the
305 probability of phosphorylation in S202 or T205. Using the program NetPhoK³⁸, we observed
306 different phosphorylation potential in S202 and T205 between G28881-G28882-G28883 (RG)
307 and A28881-A28882-C28883 (KR) (Fig. S9)

308 **CONCLUDING REMARKS:**

309 Here, we present a complete geographical and temporal worldwide distribution of SARS-CoV-2
310 haplotypes during the first five months of the pandemic. We identified 9 high-frequency
311 mutations. These important variations (asserted mainly by their frequencies) need to be tracked
312 during the pandemic.

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

317 Our geographical and temporal analysis showed that OTU_3 is currently the more frequent
318 haplotype circulating in the world. Even in North America that seems to be not the most
319 frequent in the overall analysis, seems to be that in June it is the most frequent. These results
320 highlight the importance to study mutations that characterize this haplotype, those in the
321 nucleocapsid protein R203K and G204R.

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

325 In the next months, these haplotypes description will need to be updated, identification of new
326 haplotypes could be performed by combining the identification of new frequent mutations and
327 phylogenetic analysis. We will continue monitoring the emergence of mutations that exceed our
328 proposed cut-off of 0.1-0.9 NRF_p and this information will be rapidly shared with the scientific
329 community through our web. This will also be accompanied by a continuous update of
330 haplotypes information.

331 Our weak conclusion related to age, gender, and patient status information is due to the poorly
332 organized metadata publicly available. Thus, we highlight the importance of correct
333 management and organization of genomic metadata. Regarding this, we are setting up a web
334 system where the scientific community can voluntarily share patient information associated
335 with genomic data in an organized manner. This will allow filling current gaps in the public data
336 on the correlation of haplotypes (or variants in general) with the severity of the disease, specific
337 symptoms, comorbidities, among others.

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

342 **MATERIAL AND METHODS:**

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

344 To perform the mutation frequency analysis, we first downloaded a total of 48 776 complete
345 and high coverage genomes from the GISAID database (as of July 30th, 2020). This set of genomes
346 was aligned using MAFFT with FFT-NS-2 strategy and default parameter settings³⁹. Then, we
347 removed columns that do not correspond to the region from nt 203 to nt 29674, and insertions
348 respect to the genome EPI_ISL_402125. After that, these regions were aligned using MAFFT with
349 FFT-NS-2 strategy and default parameter settings³⁹. Subalignments corresponding to genomes
350 divided by continent-month combinations was extracted and relative frequencies of each base
351 or gap in each genomic position were calculated ($RF_{p,m-c}$) using a python script. These relative
352 frequencies were multiplied by the number of cases reported in the respective continent-month
353 combination (CN_{m-c}) obtaining an estimation of the number of cases that present a virus with
354 a specific base or gap in a specific genomic position ($RF_p CN_{m-c}$). Finally, we added the
355 $RF_p CN_{m-c}$ of each subalignment and divided it by the total number of cases in the world
356 ($\sum_{m-c} RF_o CN_{m-c1} / TCN_w$). This procedure allows us to obtain a relative frequency normalized
357 by cases of each base or gap in each genomic position (NRF_p). The number of cases of each
358 country was obtained from the European Centre for Disease Prevention and Control:
359 [https://www.ecdc.europa.eu/en/publications-data/download-todays-data-geographic-](https://www.ecdc.europa.eu/en/publications-data/download-todays-data-geographic-distribution-covid-19-cases-worldwide)
360 [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 number of cases of countries with at least
361 one genome sequenced and deposited in GISAID database. Also, we just consider in the analysis
362 month-continent combinations with at least 90 genomes sequenced.

363 **Phylogenetic tree construction:**

364 Using an alignment of the 19 486 complete, high coverage genomes without ambiguities, we
365 estimated a maximum likelihood tree with IQ-TREE 2⁴⁰ using the GTR+F+R2 model of nucleotide
366 substitution^{41,42,43}, default heuristic search options, ultrafast bootstrapping with 1000
367 replicates⁴⁴ and the genome EPI_ISL_408601 as the outgroup. To generate tree figures with
368 continent or date information by tip we used the maximum likelihood tree and ggtree package
369 in R^{45,46}.

370 **OTUs determination:**

371 Positions with between 0.1 and 0.9 NRF_p were extracted from the alignment of the non-
372 ambiguities data set of 19486 and were associated with the whole-genome rooted tree using
373 the MSA function from the ggtree package^{45,46} in R. Then, we visually examined to identify the
374 major haplotypes based in these positions, designated as OTUs (Operational Taxonomic Units).
375 Haplotypes identification based in our NRF_p calculation reduced the bias of the different number
376 of genomes sequenced in each continent and each month by integrating the less biased
377 information of the number of cases. Although, other biases are more difficult, if possible, to
378 reduce or eliminate.

379 **Analysis of OTUs geographical distribution:**

380 In this analysis, we randomly separate the genomes into 6 groups of 8 129 genomes each and
381 we analyzed them independently. After that, genomes in each sample was divided by continents
382 and by months. In these divisions, OTUs relative frequencies were calculated for each OTU in
383 each month-continent combination ($O_n F_{m-c}$). Then, we multiplied these ($O_n F_{m-c}$) frequencies
384 by the number of cases corresponding to the respective month-continent (CN_{m-c}) to obtain an
385 estimation of the number of cases caused by a specific OTU in a respective month-continent
386 ($O_n CN_{m-c}$). After, these products were grouped by continents, and those from the same
387 continent were added and then divided by the total number of cases in the continent analyzed

388 $(\sum_{m-c1} O_n CN_{m-c1})/TCN_{c1}$. Thus, obtaining a frequency normalized by cases for each OTU in
389 each continent. Finally, following this procedure in each sample, we statistically compared the
390 mean of those six samples using the package “ggpubr” in R with the non-parametric Kruskal-
391 Wallis test, and pairwise statistical differences were calculated using non-parametric Wilcoxon
392 test from the same R package. The number of cases of each country was obtained from the
393 European Centre for Disease Prevention and Control:
394 [https://www.ecdc.europa.eu/en/publications-data/download-todays-data-geographic-](https://www.ecdc.europa.eu/en/publications-data/download-todays-data-geographic-distribution-covid-19-cases-worldwide)
395 [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 number of cases of countries with at least
396 one genome sequenced and deposited in GISAID database. Also, we just consider in the analysis
397 month-continent combinations with at least 90 genomes sequenced.

398 **Analysis of OTUs temporal distribution:**

399 Following a similar procedure used in the geographical analysis, we now grouped the products
400 $O_n CN_{m-c}$ by months, added them, and then divided by the total number of cases in the analyzed
401 month $(\sum_{m1-c} O_n CN_{m1-c})/TCN_{m1}$. As in the geographical analysis, the mean of the six
402 samples was statistically compared using the same procedures and with exactly the same
403 consideration of month-continent combinations.

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

405 4108 complete and high coverage genomes with patient status information were downloaded
406 from the GISAID database (as of 30th July) and classified in OTUs using python scripts. Patient
407 status information from GISAID was recategorized in four disease levels: No Informative,
408 Asymptomatic, Mild, and Severe. A table showing the GISAID patient status categorize
409 comprising our categories can be found in Figure S6.A. We calculate the relative and absolute
410 frequency of OTUs in each patient status category. Also, using all the available information on
411 gender and age in the 48 776 genomes, we calculated the relative and absolute frequency of
412 OTUs by age and gender.

413 **DATA AVAILABILITY:**

414 The data that support the findings of this study are available on request from the corresponding
415 author upon reasonable request.

416 **REFERENCES:**

- 417 1. Cuccinotta D and Vanelli M. 2020. WHO declares COVID-19 a pandemic. Vol. 91, 157-
418 160. Acta Biomedica.
- 419 2. WHO. 2020. <https://www.who.int/emergencies/diseases/novel-coronavirus-2019>
420 Retrieved on 25 August 2020.
- 421 3. Zhu N, Zhang D, Wang W, Li X, Yang B, Song J, Zhao X, Huang B, Shi W, Lu R, Niu P, Zhan
422 F, Ma X, Wang D, Xu W, Wu G, Gao G, Tan W. 2020. A novel coronavirus from patients
423 with pneumonia in China, 2019. Vol. 382(8), 727-733. The New England Journal of
424 Medicine.
- 425 4. Shu Y and McCauley J. 2017. GISAID: Global initiative on sharing all influenza data – from
426 vision to reality. Vol. 22(13), 1-3. Euro Surveillance.
- 427 5. Tang X, Wi C, Li X, Song Y, Yao X, Wu X, Duan Y, Zhang H, Wang Y, Qian Z, Cui J, Lu J.
428 2020. On the origin and continuing evolution of SARS-CoV-2. Vol. 7, 1012-1023. National
429 Science Review.

- 430 6. Rambaut A, Holmes E, Hill V, O'Toole A, Hill V, McCrone J, Ruis C, du Plessis L, Pybus O.
431 2020. A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic
432 epidemiology. 2020. Nature Microbiology. <https://doi.org/10.1038/s41564-020-0770-5>
- 433 7. Korber B, Fischer W, Gnanakaran S, Yoon H, Theiler J, Abfalterer W, Hengartner N, Giorgi
434 E, Bhattacharya T, Foley B, Hastie K, Parker M, Partridge D, Evans C, Freeman T, de Silva
435 T, McDanal C, Perez L, Tang H, Moon-Walker A, Whelan S, LaBranche C, Saphire E,
436 Montefiori D. 2020. Tracking changes in SARS-CoV-2 Spike: evidence that D614G
437 increases infectivity of the COVID-19 virus. <https://doi.org/10.1016/j.cell.2020.06.043>.
438 Cell.
- 439 8. Kern D, Sorum B, Hoel C, Sridharan S, Remis J, Toso D, Brohawn S. 2020. Cryo-EM
440 structure of the SARS-CoV-2 3a ion channel in lipid nanodiscs. bioRxiv preprint doi:
441 <https://doi.org/10.1101/2020.06.17.156554>.
- 442 9. Pachetti M, Marini B, Benedetti F, Giudici F, Mauro E, Storici P, Masciovecchio C,
443 Angeletti S, Ciccozzi M, Gallo R, Zella D, Ippodrino R. 2020. Emerging SARS-CoV-2
444 mutation hot spots include a novel RNA-dependent-RNA polymerase variant. Vol.
445 18(179), 1-9. Journal of Translational Medicine.
- 446 10. Yin C. 2020. Genotyping coronavirus SARS-CoV-2: methods and implication. Genomics.
447 <https://doi.org/10.1016/j.ygeno.2020.04.016>
- 448 11. Madhugiri R, Fricke M, Marz M, Ziebuhr J. 2014. RNA structure analysis of
449 alphacoronavirus terminal genome regions. Vol. 194, 76-89. Virus Research.
- 450 12. Masters P. 2019. Coronavirus genomic RNA packaging. Vol. 537, 198-207. Virology.
- 451 13. Mukherjee M and Goswami S. 2020. Global cataloguing of variations in untranslated
452 regions of viral genome and prediction of key host RNA binding protein-microRNA
453 interactions modulating genome stability in SARS-CoV-2. bioRxiv preprint doi:
454 <https://doi.org/10.1101/2020.06.09.134585>
- 455 14. Kuo P, Chiang C, Wang Y, Doudeva L, Yuan H. 2014. The crystal structure of TDP-43
456 RRM1-DNA complex reveals the specific recognition for UG- and TG-rich nucleic acids.
457 Vol. 42(7), 4712-4722. Nucleic Acids Research.
- 458 15. Lee E, Lee V, Trojanowski J. 2011. Gains or losses: molecular mechanisms of TDP43-
459 mediated neurodegeneration. Vol. 13(1), 38-50. Nature Reviews Neuroscience.
- 460 16. Graham R, Sims A, Brockway S, Baric S, Denison M. 2005. The nsp2 replicase proteins of
461 murine hepatitis virus and severe acute respirator syndrome coronavirus are
462 dispensable for viral replication. Vol. 79(21), 13399-13411. Journal of Virology.
- 463 17. Cornillez-Ty C, Liao L, Yates J, Kuhn P, Buchmeier M. 2009. Severe acute respiratory
464 syndrome coronavirus nonstructural protein 2 interacts with a host protein complex
465 involved in mitochondrial biogenesis and intracellular signaling. Vol. 83(19), 10314-
466 10318. Journal of Virology.
- 467 18. Wang S, Nath N, Adlam M, Chellappan S. 1999. Prohibitin, a potential tumor suppressor,
468 interacts with RB and regulates E2F function. Vol. 18, 3501-3510. Oncogene.
- 469 19. Rajalingam K, Wunder C, Brinkmann V, Churin Y, Hekman M, Sievers C, Rapp U, Rudel T.
470 2005. Prohibitin is required for RAS-induced RAF-MEK-ERK activation and epithelial cell
471 migration. Vol. 7(8), 837-843. Nature Cell Biology.
- 472 20. Sun L, Liu L, Yang X, Wu Z. 2004. Akt binds prohibitin 2 and relieves its repression of
473 MyoD and muscle differentiation. Vol. 117(14), 3021-3029. Journal of Cell Science
- 474 21. Fusaro G, Dasgupta P, Rastogi S, Joshi B, Chellappan S. 2003. Prohibitin induces the
475 transcriptional activity of p53 and is exported from the nucleus upon apoptotic signaling.
476 Vol. 278(48), 47853-47861. The Journal of Biological Chemistry.

- 477 22. Merkwirth C and Langer T. 2008. Prohibitin function within mitochondria: essential roles
478 for cell proliferation and cristae morphogenesis. Vol. 1793, 27-32. *Biochimica et*
479 *Biophysica Acta*.
- 480 23. Mauro V and Chapel S. 2014. A critical analysis of codon optimization in human
481 therapeutics. Vol. 20(11), 604-613. *Trends in Molecular Medicine*.
- 482 24. Gu H, Chu D Peiris M, Poon L. 2020. Multivariate Analyses of Codon Usage of SARS-CoV-
483 2 and other betacoronaviruses. bioRxiv preprint doi:
484 <https://doi.org/10.1101/2020.02.15.950568>.
- 485 25. Chand G and Azad G. 2020. Identification of novel mutations in RNA-dependent RNA
486 polymerases of SARS-CoV-2 and their implications. bioRxiv preprint doi:
487 <https://doi.org/10.1101/2020.05.05.079939>.
- 488 26. Zhang L, Jackson C, Mou H, Ojha A, Rangarajan E, Izard T, Farzan M, Choe H. 2020. The
489 D614G mutation in the SARS-CoV-2 spike protein reduces S1 shedding and increases
490 infectivity. bioRxiv preprint doi: <https://doi.org/10.1101/2020.06.12.148726>.
- 491 27. Castaño-Rodríguez C, Honrubia J, Gutierrez-Alvarez J, DeDiego M, Nieto-Torres J,
492 Jimenez-Guardeño J,, Regla-Nava J, Fernandez-Delgado R, Verdia-Báguena C, Queralt-
493 Martín M, Kochan G, Perlman S, Aguilella V, Sola I, Enjuanes L. 2018. Role of severe acute
494 respiratory syndrome coronavirus viroporins E, 3a, and 8a in replication and
495 pathogenesis. Vol. 9(3), 1-23. *American Society for Microbiology*.
- 496 28. Siu K, Yuen K, Castaño-Rodríguez C, Ye Z, Yeung M, Fung S, Yuan S, Chan C, Yuen K,
497 Enjuanes L, Jin D. 2019. Severe acute respiratory syndrome coronavirus ORF3a protein
498 activates the NLRP3 inflammasome by promoting TRAF3-dependent ubiquitination of
499 ASC. Vol. 33. 8865-8877. *The FASEB Journal*.
- 500 29. Chan C, Tsoi H, Chan W, Zhai S, Wong C, Yao X, Chan W, Tsui S, Chan H. 2009. The ion
501 channel activity of the SARS-coronavirus 3a protein is linked to its pro-apoptotic
502 function. Vol. 41, 2232-2239. *The International Journal of Biochemistry and Cell Biology*
- 503 30. Yue Y, Nabar N, Shi C, Kamenyeva O, Xiao X, Hwang I, Wang M, Kehrl J. 2018. SARS-
504 Coronavirus open reading frame-3a drives multimodal necrotic cell death. Vol. 9, 1-15.
505 *Cell Death and Disease*.
- 506 31. Padhan K, Tanwar C, Hussain A, Hui P, Lee M, Cheung C, Malik J, Jameel S. 2007. Severe
507 acute respiratory syndrome coronavirus Orf3a protein interacts with caveolin. Vol. 88,
508 3067-3077. *Journal of General Virology*.
- 509 32. Griffiths G and Simons K. 1986. The trans Golgi network: sorting at the exit site of the
510 golgi complex. Vol. 234, 438-443. *Science*.
- 511 33. Lau S, Fen Y, Cheng H, Luk H, Yang W, Li K, Zhang Y, Huang Y, Song Z, Chow W, Fan R,
512 Ahmed S, Yeung H, Lam C, Cai J, Wong Chan J, Yuen K, Zhang H, Woo P. 2015. Severe
513 acute respiratory syndrome (SARS) coronavirus ORF8 protein is acquired from SARS-
514 related coronavirus from greater horseshoe bats through recombination. Vol. 89(20),
515 10532-10547. *Journal of Virology*.
- 516 34. Zhang Y, Zhang J, Chen Y, Luo B, Yuan Y, Huang F, Yang T, Yu F, Liu J, Liu B, Song Z, Chen
517 J, Pan T, Zhang X, Li Y, Li R, Huang W, Xiao F, Zhang H. 2020. The Orf8 protein of SARS-
518 CoV-2 mediates immune evasion through potentially downregulating MHC-1. bioRxiv
519 preprint doi: <https://doi.org/10.1101/2020.05.24.111823>.
- 520 35. Chang C, Hsu Y, Chang Y, Chao F, Wu M, Huang Y, Hu C, Huang T. 2009. Multiple nucleic
521 acid binding sites and intrinsic disorder of severe acute respiratory syndrome
522 coronavirus nucleocapsid protein implications for ribonucleocapsid protein packaging.
523 Vol. 83(5), 2255-2264. *Journal of Virology*.

- 524 36. Luo H, Chen Q, Chen J, Chen K, Shen X, Jiang H. 2005. The nucleocapsid protein of SARS
525 coronavirus has a high binding affinity to the human cellular heterogeneous nuclear
526 ribonucleoprotein A1. Vol. 579, 2623-2628. FEBS letters.
- 527 37. Chang C, Chen C, Chiang M, Hsu Y, Huang T. 2013. Transient oligomerization of the SARS-
528 CoV N protein – Implication for virus ribonucleoprotein packaging. Vol. 8(5), e65045.
529 PlosONE.
- 530 38. Blom N, Sicheritz-Pontén T, Gupta R, Gammeltoft S, Brunak S. 2004. Prediction of post-
531 translational glycosylation and phosphorylation of proteins from the amino acid
532 sequence. Vol. 4, 1633-1649. Proteomics.
- 533 39. Katoh K, Misawa K, Kuma K, Miyata T. 2002. MAFFT: a novel method for rapid multiple
534 sequence alignment based on fast Fourier transform. Vol. 30(14), 3059-3066. Nucleic
535 Acids Research.
- 536 40. Minh B, Schmidt H, Chernomor O, Schrempf D, Woodhams M, von Haeseler A, Lanfear
537 R. 2020. IQ-TREE 2: New models and efficient methods for phylogenetic inference in the
538 genomic era. Vol. 37, 1530-1534. Molecular Biology and Evolution.
- 539 41. Tavaré S. 1986. Some mathematical questions in biology: DNA sequence analysis.
540 Lectures on mathematics in the life sciences. Vol. 17, 57-86. American Mathematical
541 society.
- 542 42. Soubrier J, Steel M, Lee M, Sarkissian C, Guindon S, Ho S, Cooper A. 2012. The influence
543 of rate heterogeneity among sites on the time dependence of molecular rates. Vol.
544 29(11), 3345-3358. Molecular Biology and Evolution.
- 545 43. Yang Z. 1995. A space-time process model for the evolution of DNA sequences. Vol. 139,
546 993-1005. Genetics
- 547 44. Hoang D, Chernomor O, Haeseler A, Minh B, Vinh L. 2017. UFBoot2: Improving the
548 ultrafast bootstrap approximation. Vol. 35(2), 518-522. Molecular Biology and
549 Evolution.
- 550 45. Yu G. 2020. Using ggtree to visualize data on tree-like structures. Vol. 69, 1-18- Current
551 Protocols in Bioinformatics.
- 552 46. Yu G, Smith D, Zhu H, Guan Y, Lam T. 2017. GGTREE: an R package for visualization and
553 annotation of phylogenetic trees with their covariates and other associated data. Vol. 8,
554 28-36. Methods in Ecology and Evolution.

555 **Competing interests:**

556 The authors declare no competing interests

557 **Acknowledgements:**

558 We thank Professor Shaker Chuck Farah (Institute of Chemistry – University of Sao Paulo) for
559 English writing corrections and helpful comments. Also, we thank Professors Aline Maria da Silva
560 (Institute of Chemistry – University of Sao Paulo), Joao Renato Rebelo Pinho (Albert Einstein
561 Hospital – Sao Paulo) and PhD(c). Deyvid Amgarten (Albert Einstein Hospital – Sao Paulo) for its
562 helpful comments. To the Ricardo Palma University High-Performance Computational Cluster
563 (URPHPC) managers Gustavo Adolfo Abarca Valdiviezo and Roxana Paola Mier Hermoza at the
564 Ricardo Palma Informatic Department (OFICIC) for their contribution in programs and remote
565 use configuration of URPHPC. To Gladys Arevalo Chong for her figure style suggestions.

566

567

568 **Figure 1. Five haplotypes (or OTUs) based in nine positions can classify 96.5 % of the**
569 **genomes.** A) Table showing haplotype of each OTU, regions, and aminoacids changes caused
570 by these mutations. B) Rooted tree of 19486 SARS-CoV-2 complete and non-ambiguous
571 genomes associated with an alignment of nine genomic positions (241, 1059, 3037, 14408,
572 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).

573 **Figure 2. By cases normalized continent distribution of OTUs during seven months of the**
574 **pandemic.** A) Unrooted tree of complete non-ambiguous genomes, tips were colored
according to OTUs, and points in each tip were colored according to the continent. B-G)
575 Boxplots of normalized relative frequencies of OTUs in each continent (B, North America; C,
576 South America; D, Europe; E, Asia; F, Oceania; G, Africa).

577 **Figure 3. By cases normalized temporal distribution of OTUs showed OTU_3 as the**
578 **currently most prevalent.** A) Rooted tree of complete non-ambiguous genomes showing
temporal distribution. Tips were colored by OTUs and points in each tip were colored
579 according to the isolation date. B-E) Boxplot of OTUs distribution in each month (B, February;
580 C, March; D, April; E, May; F, June; G, July).

581 **Figure 4. Age, gender, and patient status distribution by OTUs do not show preferences for**
582 **patient characteristics.** A) Relative and absolute frequencies of OTUs distribution by age. B)
Age distribution was grouped by ranges and relative and absolutes frequency by OTUs is
583 showed. C) OTUs distribution by gender. D) Relative and absolute frequencies of OTUs by
patient status categories.

584 **Figure 5. P323L could impact the stability of Nsp12 without disturbing its overall structure.**
A) Structure of RNA-dependent RNA polymerase complex (PDB ID: 6YYT). Chains (Nsp12,
585 Nsp7, Nsp8, RNA) are distinguished by colors. Helix 10, Beta-sheet 3, Turn 10-3, and P323
586 also are differentially colored. B) Structure in A rotated 90 degrees. C) Zoom of the red box
in B showed P322 and P323 in the center of Turn 10-3. D) Turn 10-3 with side chains of P323,
587 L324, and F396 in sphere representation to highlight the distance between side chains of
P323 and L324. E) P323 in D was computationally replaced by L323. Now, distances between
588 the methyl group of leucine are shorter with L323.

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

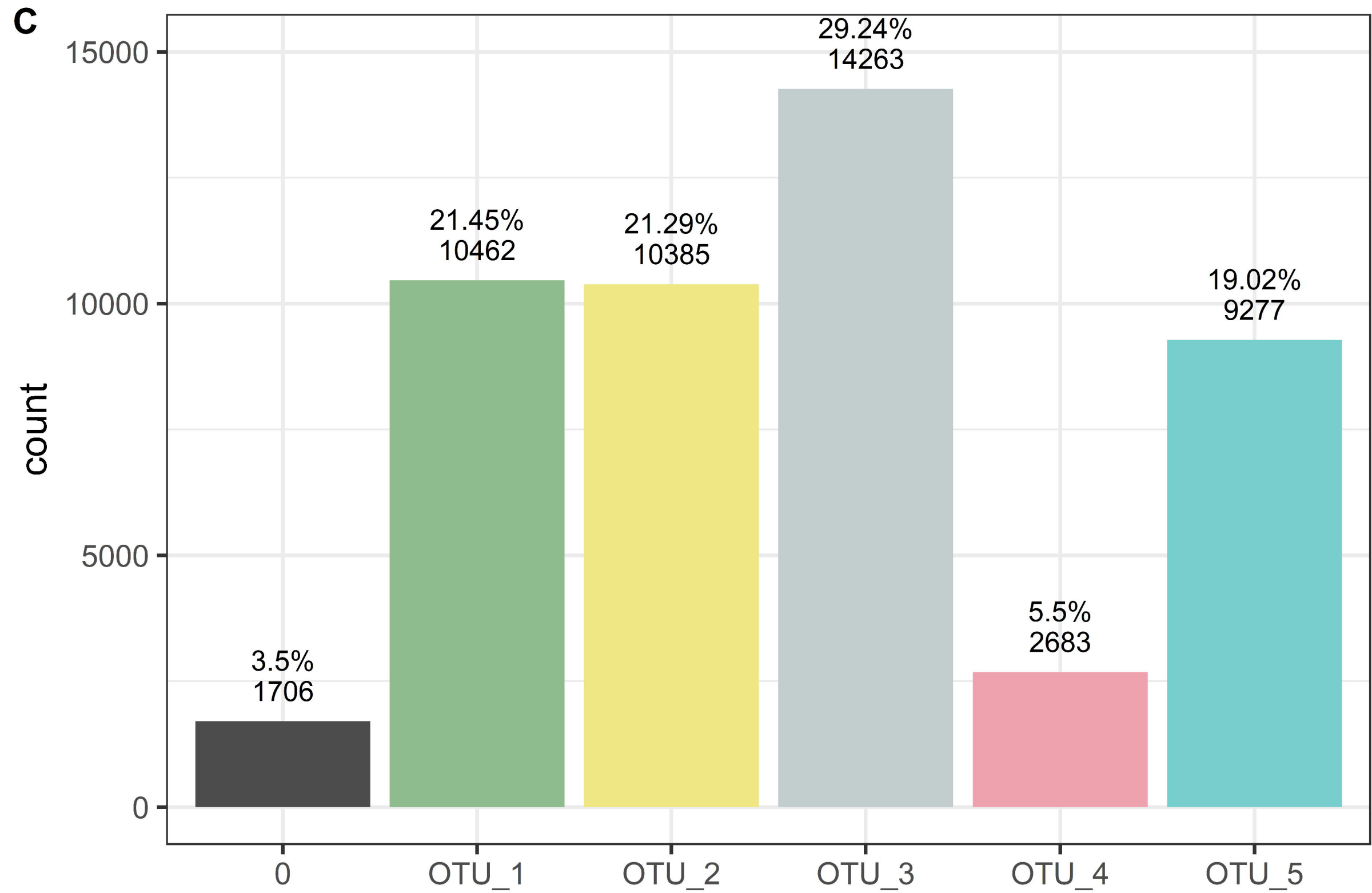
594 **Figure 7. Orf3a Q57H does not modify pore constriction distances but electrostatics**
595 **distribution.** A) Structure of the Orf3a dimer (PDB ID: 6XDC) colored by domains. The right
of A shows the same structure but in an upper view. B) Orf3a showing the central pore, in
596 the red box the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

A

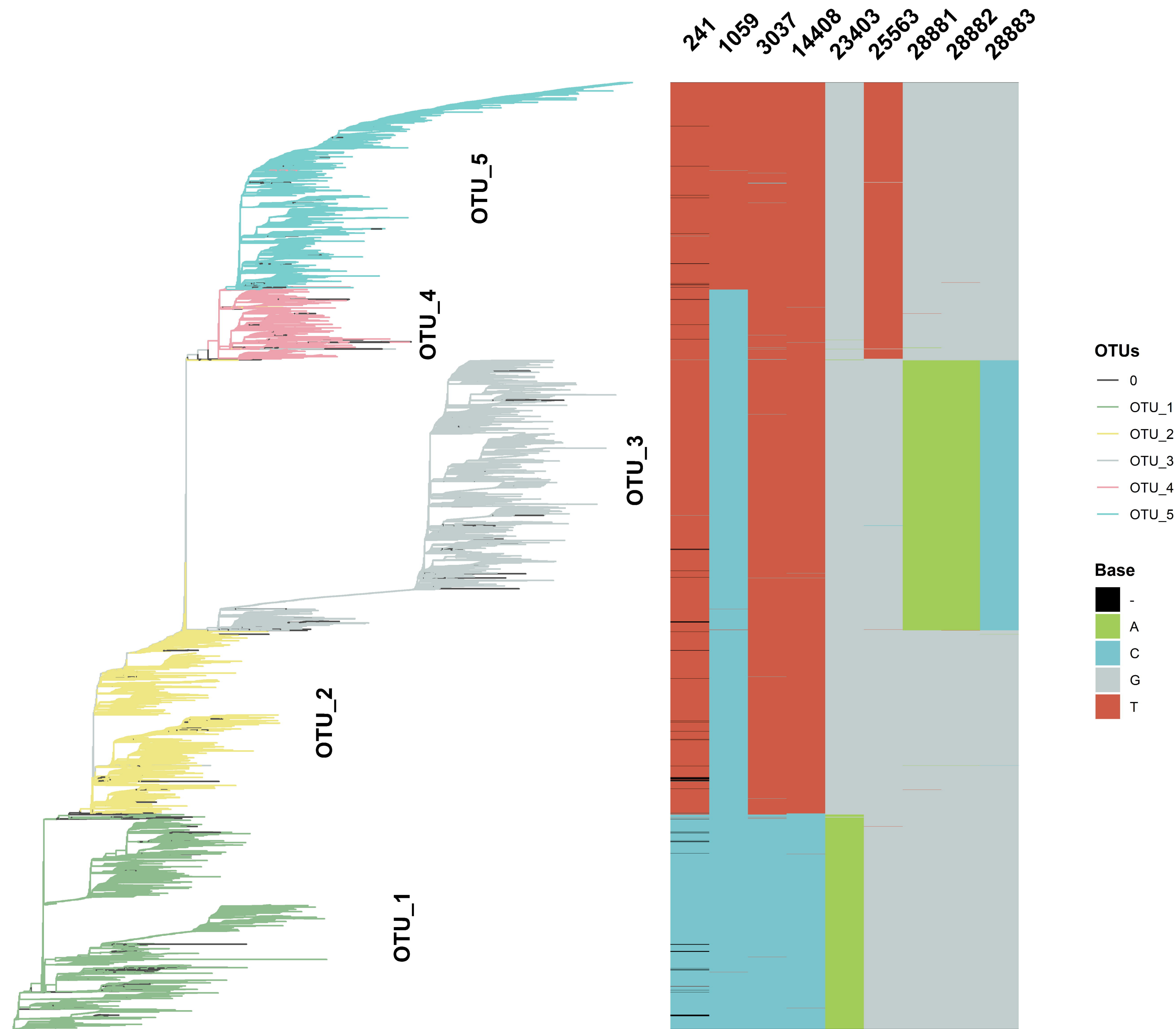
Position	OTU_1	OTU_2	OTU_3	OTU_4	OTU_5	Region	AA_change
241	C	T	T	T	T	5'UTR	''
1059	C	C	C	C	T	Nsp2	T85I
3037	C	T	T	T	T	Nsp3	Syn
14408	C	T	T	T	T	Nsp12	P323L
23403	A	G	G	G	G	S	D614G
25563	C	G	G	T	T	Orf3a	Q57H
28881	G	G	A	G	G	N	R203K
28882	G	G	A	G	G	N	R203K
28883	G	G	C	G	G	N	G204R

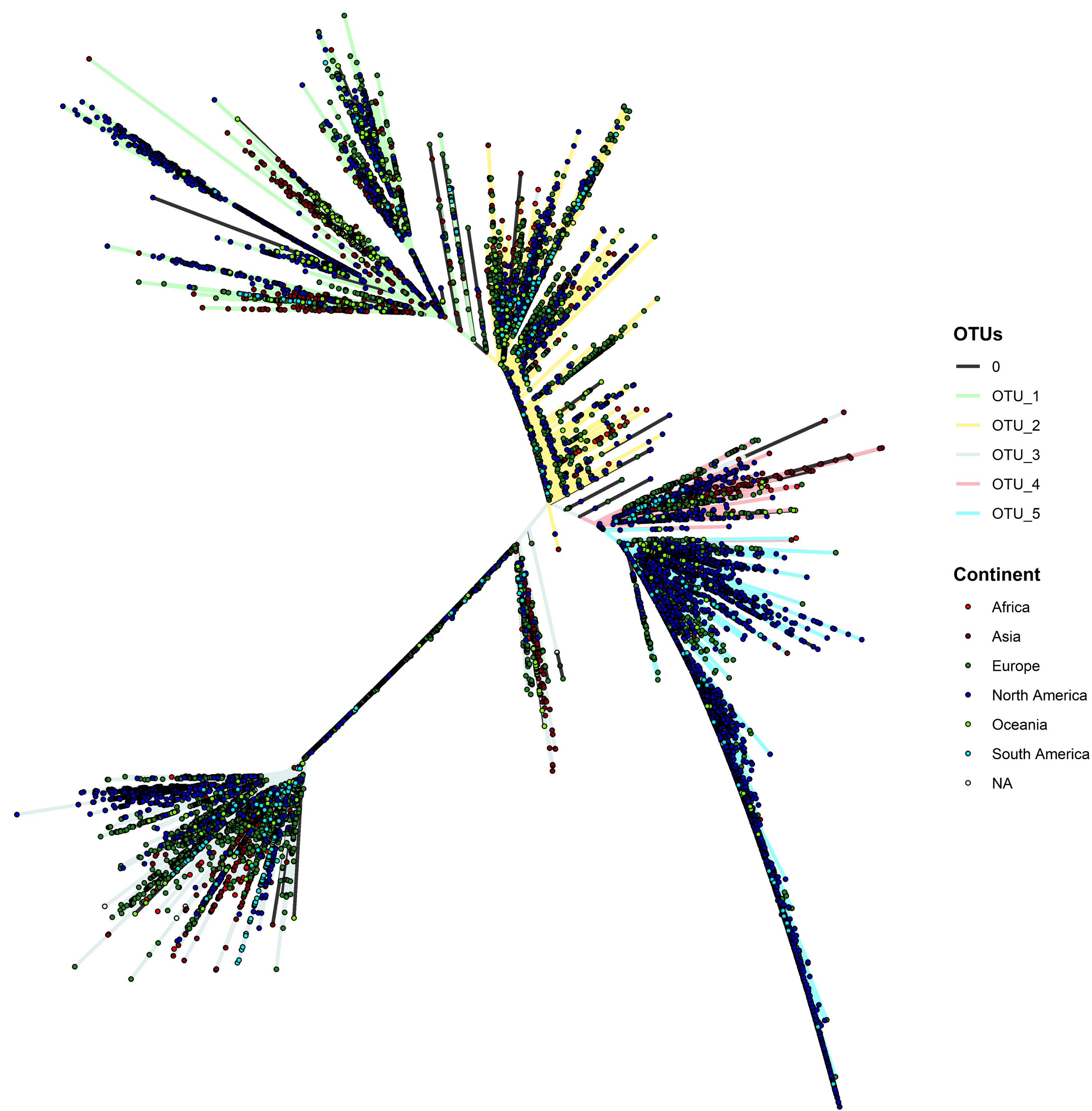
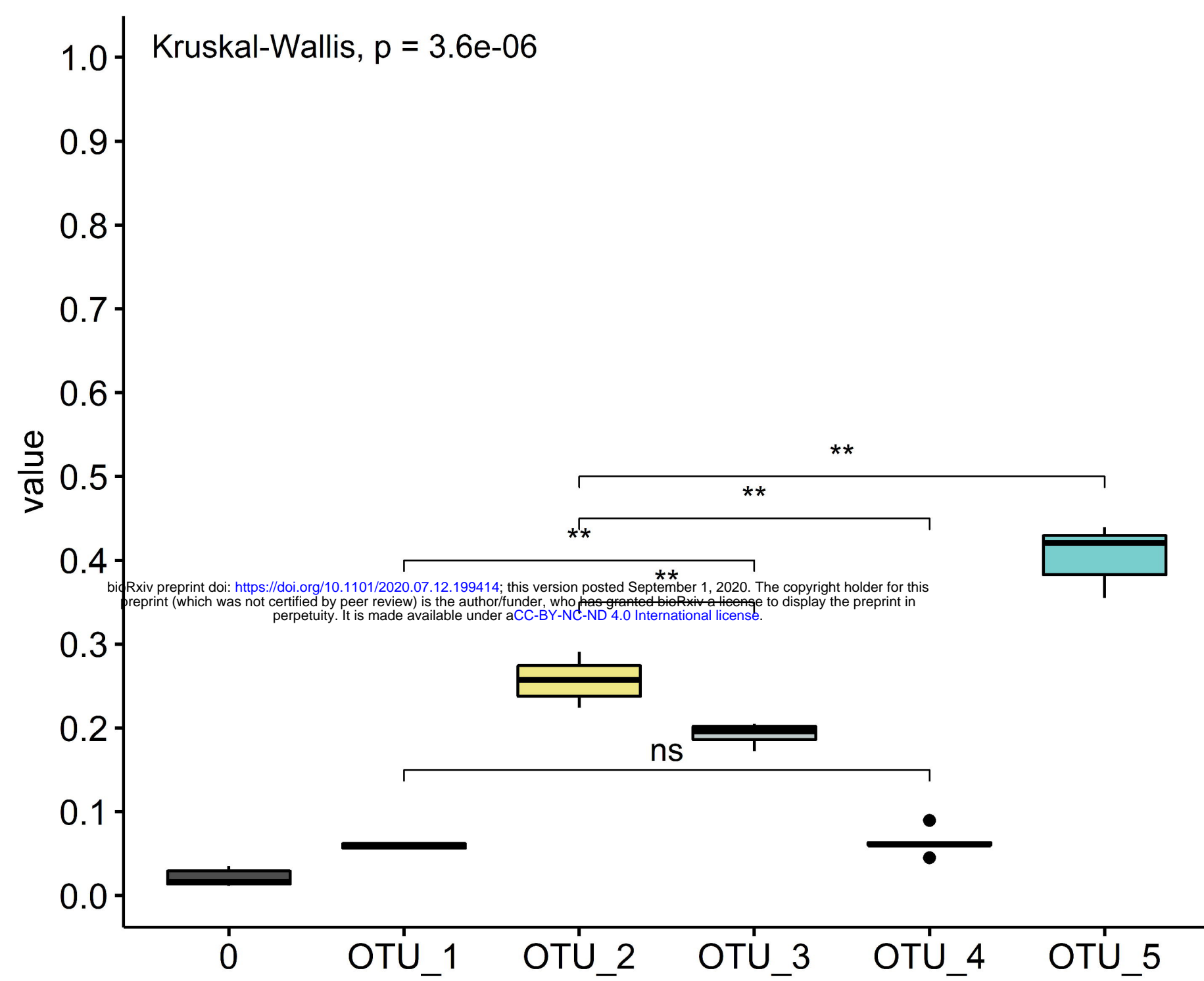
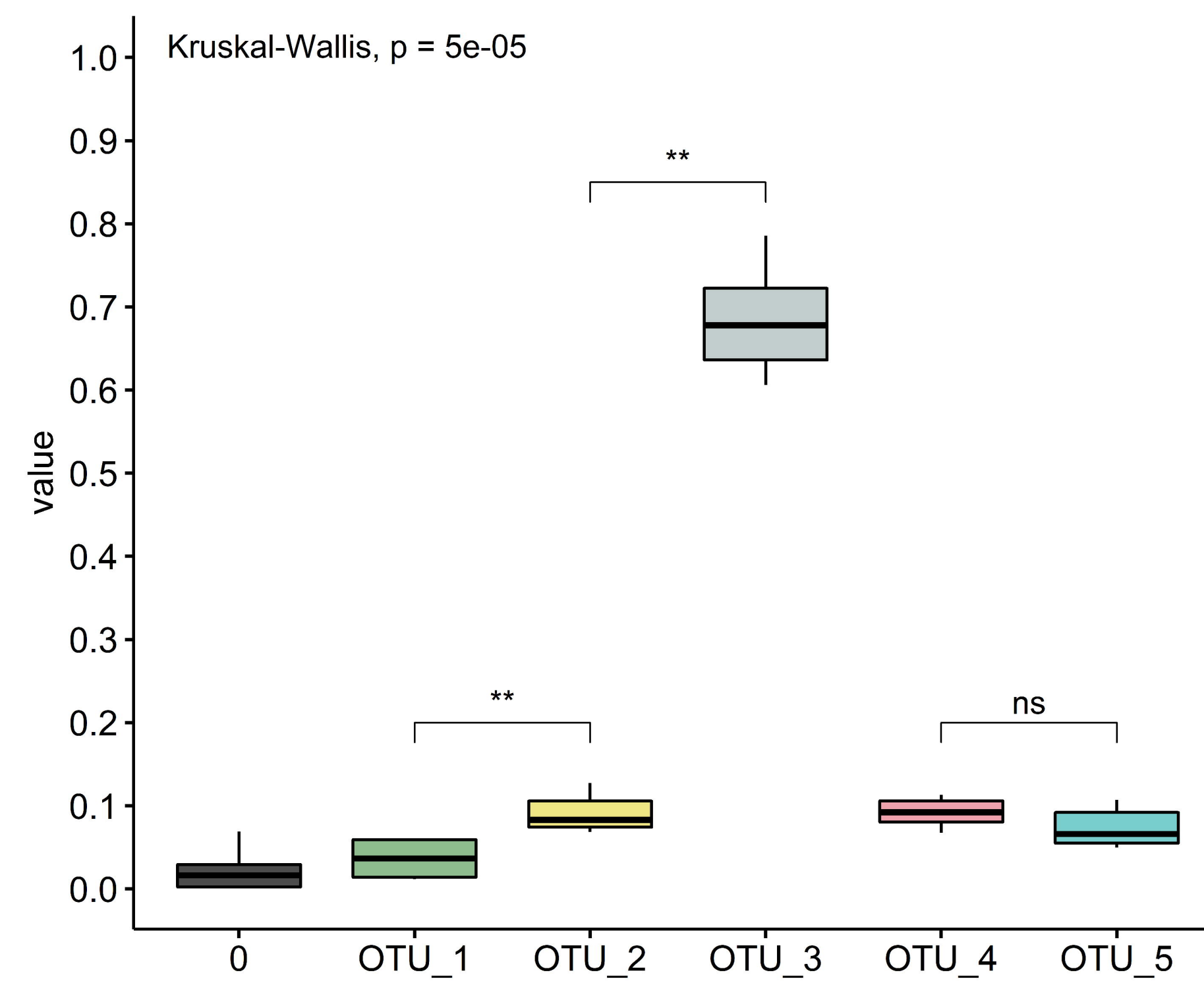
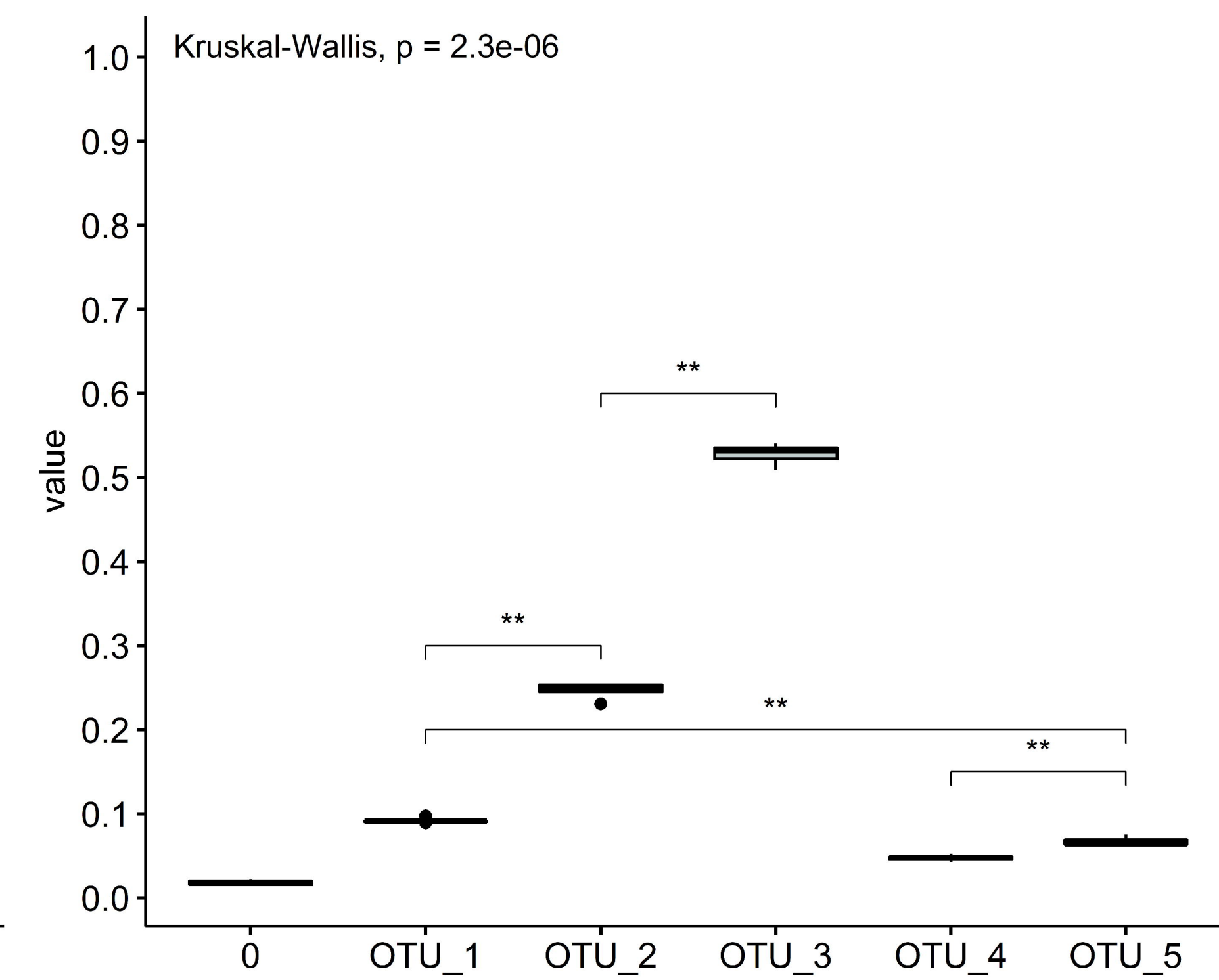
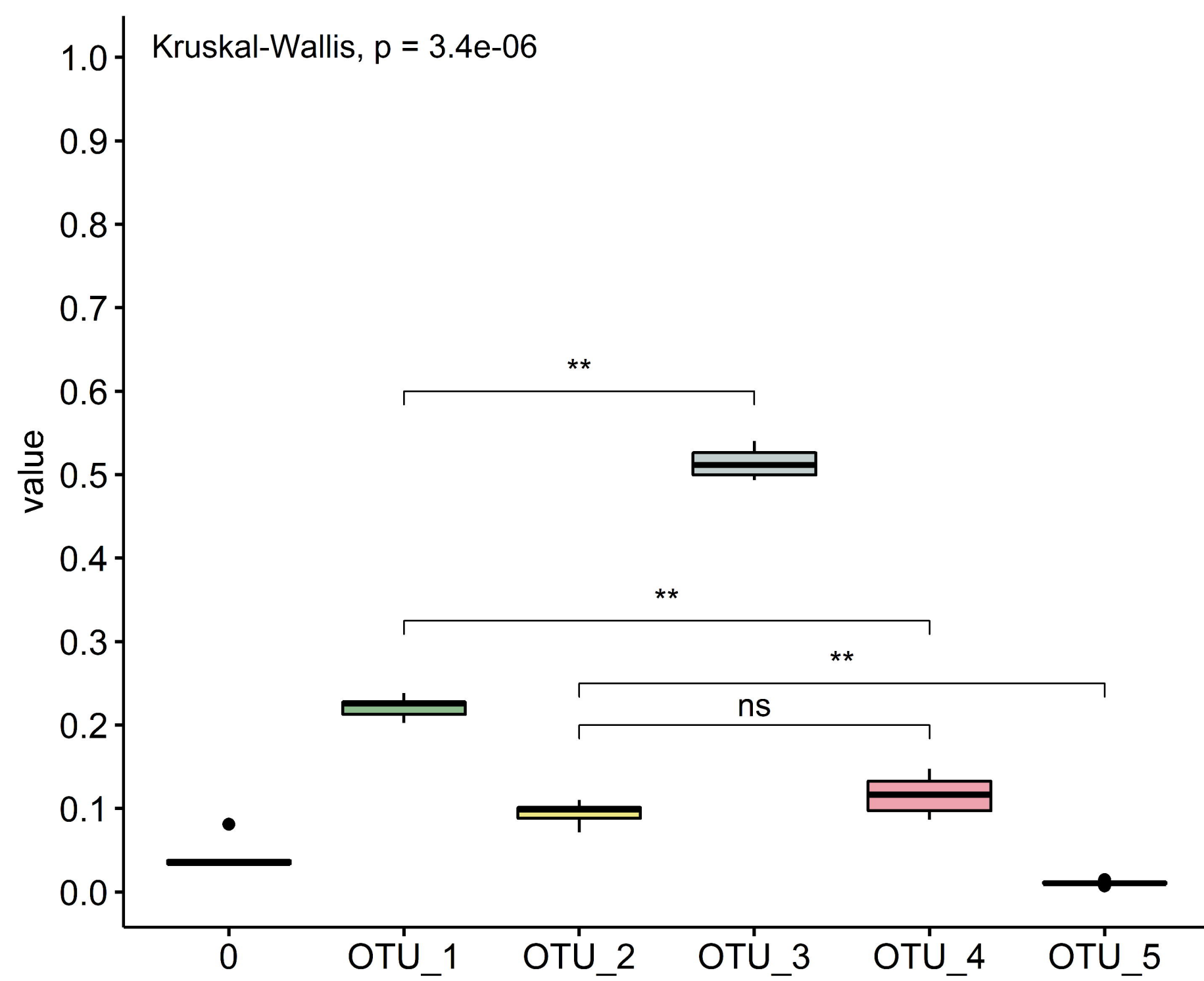
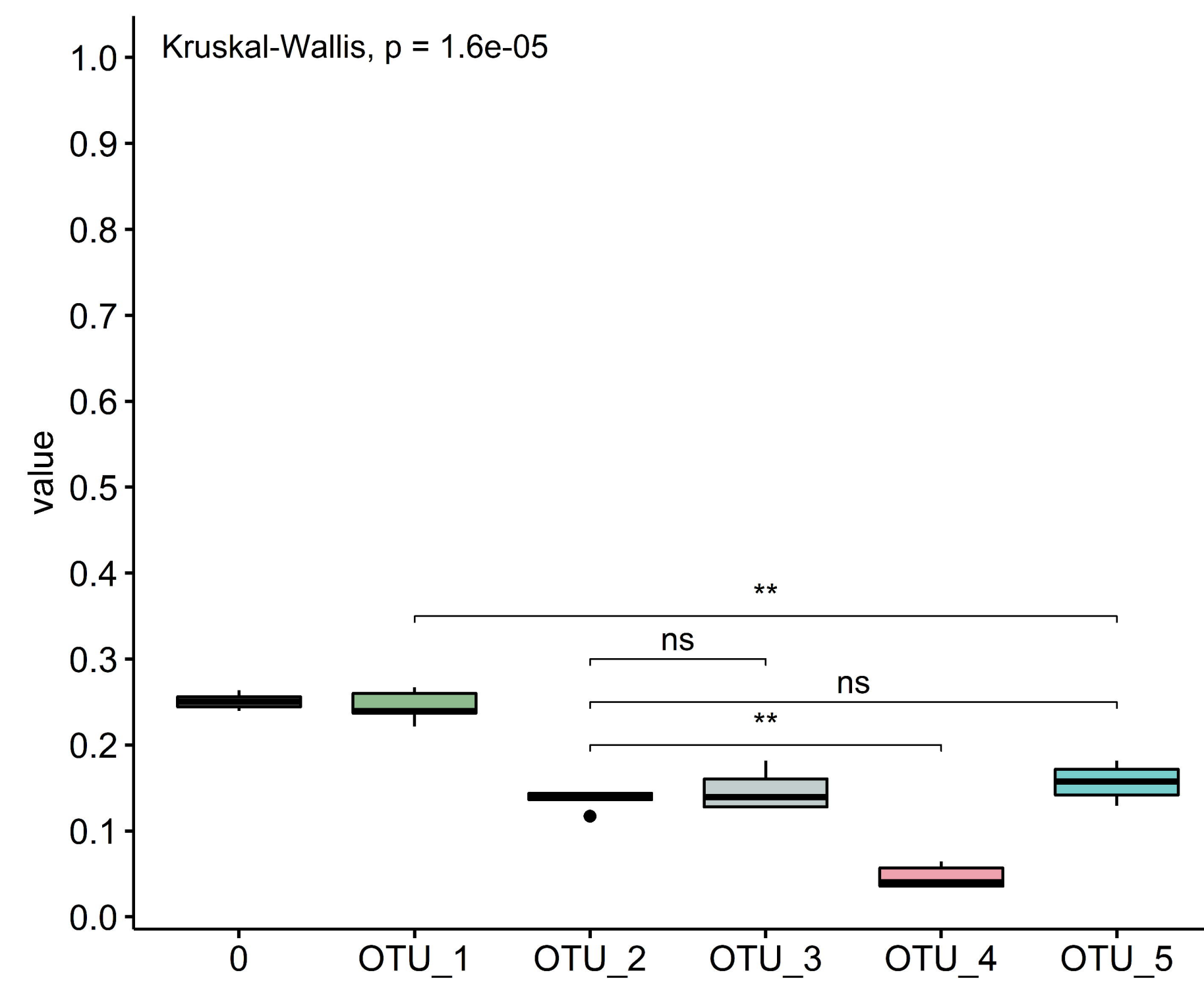
bioRxiv preprint doi: <https://doi.org/10.1101/2020.07.12.199414>; this version posted September 1, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

C



B



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