



Global Geographic and Temporal Analysis of SARS-CoV-2 Haplotypes Normalized by COVID-19 Cases During the Pandemic

Santiago Justo Arevalo^{1,2*}, Daniela Zapata Sifuentes¹, César J. Huallpa³, Gianfranco Landa Bianchi¹, Adriana Castillo Chávez¹, Romina Garavito-Salini Casas¹, Guillermo Uceda-Campos⁴ and Roberto Pineda Chavarria¹

OPEN ACCESS

Edited by:

Basil Britto Xavier,
University of Antwerp, Belgium

Reviewed by:

Ritesh Tandon,
University of Mississippi,
United States
Ludmila Chistoserdova,
University of Washington,
United States

*Correspondence:

Santiago Justo Arevalo
santiago.jus.are@usp.br;
santiago.justo@urp.edu.pe

Specialty section:

This article was submitted to
Evolutionary and Genomic
Microbiology,
a section of the journal
Frontiers in Microbiology

Received: 30 September 2020

Accepted: 25 January 2021

Published: 17 February 2021

Citation:

Justo Arevalo S, Zapata D, Huallpa CJ, Landa G, Castillo Chávez A, Garavito-Salini R, Uceda-Campos G and Pineda RC (2021) Global Geographic and Temporal Analysis of SARS-CoV-2 Haplotypes Normalized by COVID-19 Cases During the Pandemic. *Front. Microbiol.* 12:612432. doi: 10.3389/fmicb.2021.612432

¹Facultad de Ciencias Biológicas, Universidad Ricardo Palma, Lima, Peru, ²Department of Biochemistry, Institute of Chemistry, University of São Paulo, São Paulo, Brazil, ³Facultad de Ciencias, Universidad Nacional Agraria La Molina, Lima, Peru, ⁴Facultad de Ciencias Biológicas, Universidad Nacional Pedro Ruiz Gallo, Lambayeque, Peru

Since the identification of SARS-CoV-2, a large number of genomes have been sequenced with unprecedented speed around the world. This marks a unique opportunity to analyze virus spreading and evolution in a worldwide context. Currently, there is not a useful haplotype description to help to track important and globally scattered mutations. Also, differences in the number of sequenced genomes between countries and/or months make it difficult to identify the emergence of haplotypes in regions where few genomes are sequenced but a large number of cases are reported. We propose an approach based on the normalization by COVID-19 cases of relative frequencies of mutations using all the available data to identify major haplotypes. Furthermore, we can use a similar normalization approach to tracking the temporal and geographic distribution of haplotypes in the world. Using 171,461 genomes, we identify five major haplotypes or operational taxonomic units (OTUs) based on nine high-frequency mutations. OTU_3 characterized by mutations R203K and G204R is currently the most frequent haplotype circulating in four of the six continents analyzed (South America, North America, Europe, Asia, Africa, and Oceania). On the other hand, during almost all months analyzed, OTU_5 characterized by the mutation T85I in nsp2 is the most frequent in North America. Recently (since September), OTU_2 has been established as the most frequent in Europe. OTU_1, the ancestor haplotype, is near to extinction showed by its low number of isolations since May. Also, we analyzed whether age, gender, or patient status is more related to a specific OTU. We did not find OTU's preference for any age group, gender, or patient status. Finally, we discuss structural and functional hypotheses in the most frequently identified mutations, none of those mutations show a clear effect on the transmissibility or pathogenicity.

Keywords: SARS-CoV-2, COVID-19, viral pandemic, phylogenomic, global analysis, epidemiology, haplotypes, operational taxonomic units

INTRODUCTION

COVID-19 was declared a pandemic by the World Health Organization on March 11th, 2020 (Cucinotta and Vanelli, 2020), with around 71 million cases and 1.6 million deaths around the world (December 14th, 2020; WHO, 2020), quickly becoming the most important health concern in the world. Several efforts to produce vaccines, drugs, and diagnostic tests to help in the fight against SARS-CoV-2 are being mounted in a large number of laboratories all around the world.

Since the publication on January 24th, 2020 of the first complete genome sequence of SARS-CoV-2 from China (Zhu et al., 2020), thousands of genomes have been sequenced in a great number of countries on all six continents and were made available in several databases. This marks a milestone in scientific history and gives us an unprecedented opportunity to study how a specific virus evolves in a worldwide context. As of November 30th, 2020, the global initiative on sharing all influenza data (GISAID) database (Shu and McCauley, 2017) contained 171,461 genomes with at least 29,000 sequenced bases.

Several analyses have been performed to identify SARS-CoV-2 variants around the world, most of them on a particular group of genomes using limited datasets (For example, Castillo et al., 2020; Franco-Muñoz et al., 2020; Maitra et al., 2020; Saha et al., 2020). In March 2020, two major lineages were proposed based on position 8,782 and 28,144 using a data set of 103 genomes (Tang et al., 2020) which was followed by a particularly interesting proposal that identified the same major lineages (named A and B) and other sublineages (Rambaut et al., 2020).

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

To perform a reasonable analysis of the worldwide temporal and geographical distribution of SARS-CoV-2 haplotypes, we need to take into account the differences in the number of sequenced genomes in months and countries or continents. Thus, we first used a data set of 171,461 complete genomes to estimate the worldwide relative frequency of nucleotides in each SARS-CoV-2 genomic position and found nine mutations with respect to the reference genome EPI_ISL_402125 with normalized relative frequencies (NRFp) representing to be present in more than 9,500,000 COVID-19 cases. After that, using a total of 109,953 complete genomes without ambiguous nucleotides from GISAID, we performed a phylogenetic analysis and correlated the major branches with SARS-CoV-2 variants which can be classified into five haplotypes or operational taxonomic units (OTUs) based on the distribution of the nine identified nucleotide positions in our NRFp analysis. After that, we analyzed the geographical and temporal worldwide distribution of OTUs normalized by the number of COVID-19 cases. Also, we attempt to correlate these OTUs with patient status, age, and gender information. Finally, we discuss the current hypothesis of the most frequent mutations on protein

structure and function. All this information will be continuously updated in our publicly available web-page.¹

MATERIALS AND METHODS

Normalized Frequency Analysis of Each Base or Gap by Genomic Position

To perform the mutation frequency analysis, we first downloaded a total of 171,461 complete and high coverage genomes from the GISAID database (as of November 30th, 2020). This set of genomes was aligned using ViralMSA using default parameter settings, and EPI_ISL_402125 SARS-CoV-2 genome from nt 203 to nt 29,674 as the reference sequence (Li, 2018; Moshiri, 2020). Subalignments corresponding to genomes divided by continent-month combinations were extracted and relative frequencies of each base or gap in each genomic position were calculated ($RF_{p,m-c}$) using a python script. These relative frequencies were multiplied by the number of cases reported in the respective continent-month combination (CN_{m-c}) obtaining an estimation of the number of cases that present a virus with a specific base or gap in a specific genomic position ($RF_p CN_{m-c}$). Finally, we added the $RF_p CN_{m-c}$ of each subalignment and divided it by the total number of cases

in the world $\left(\sum_{m-c} RF_p CN_{(m-c)} \right) / TCN_w$. This procedure allows

us to obtain a relative frequency normalized by cases of each base or gap in each genomic position (NRF_p). The number of cases of each country was obtained from the European Centre for Disease Prevention and Control: <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 one genome sequenced and deposited in GISAID database. Also, we just consider in the analysis month-continent combinations with at least 90 genomes sequenced.

Phylogenetic Tree Construction

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

OTUs Determination

Mutations respect to EPI_ISL_402125 with NRFp greater than 0.18 were extracted from the alignment of the non-ambiguous data set of 109,953 genomes and were associated with the

¹<http://sarscov2haplofinder.urp.edu.pe/>

whole-genome rooted tree using the MSA function from the `ggtree` package (Yu et al., 2017; Yu, 2020) in R. Then, we visually examined to identify the major haplotypes based on these positions, designated as OTUs. Haplotypes identification based on our NRFp calculation reduced the bias of the different number of genomes sequenced in each continent and each month by integrating the less biased information of the number of cases. Although, other biases are more difficult, if possible, to reduce or eliminate.

Analysis of OTUs Geographical Distribution

In this analysis, we randomly separate the genomes into six samples of 28,576 genomes each. Genomes in each sample were divided by continents and by months. In these divisions, OTUs relative frequencies were calculated for each OTU in each month-continent combination ($O_n F_{m-c}$). Then, we multiplied these ($O_n F_{m-c}$) frequencies by the number of cases corresponding to the respective month-continent (CN_{m-c}) to obtain an estimation of the number of cases caused by a specific OTU in a respective month-continent ($O_n CN_{m-c}$). After, these products were grouped by continents, and those from the same continent were added and then divided by the total number of cases in the continent analyzed

$$\left(\sum_{m-c} O_n CN_{m-c} \right) / TCN_{c_i} .$$

Thus, obtaining a frequency normalized

by cases for each OTU in each continent. Finally, following this procedure in each sample, we statistically compared the mean of those six samples using the package “`ggpubr`” in R with the non-parametric Kruskal-Wallis test, and pairwise statistical differences were calculated using non-parametric Wilcoxon test from the same R package. The number of cases of each country was obtained from the European Centre for Disease Prevention and Control: <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 one genome sequenced and deposited in GISAID database. Also, we just consider in the analysis month-continent combinations with at least 90 genomes sequenced.

Analysis of OTUs Temporal Distribution

Following a similar procedure used in the geographical analysis, we now grouped the products $O_n CN_{m-c}$ by months, added them, and then divided by the total number of cases in the analyzed month

$$\left(\sum_{m-c} O_n CN_{m-c} \right) / TCN_{m_i} .$$

As in the geographical analysis, the mean of the six samples was statistically compared using the same procedures and with exactly the same considerations of month-continent combinations.

Analysis of Age, Gender, and Patient Status With OTUs Distribution

We determine if OTUs have a preference for age or gender, or cause a COVID-19 with a specific severity. For patient status and age information, we selected populations with at least 45

genomes in the category to analyze and at least two times the total number of genomes (for example, Asia – February has 58 asymptomatic genomes and 613 total genomes). For the gender analysis, we selected sample populations with at least 250 genomes in the category to analyze and at least two times the total number of genomes (for example, USA – March has 2,079 genomes from female patients and 9,287 genomes with or without gender information). In each selected sample, we used the total data (all genomes corresponding to that continent-month combination) and the data with category information (for example, male, female, asymptomatic, severe, 16–30 years, etc.). We randomly divided these two groups of genomes into three samples and calculated OTUs frequencies. The mean of the frequency of each OTU was compared between the two groups using the non-parametric Wilcoxon or Kruskal-Wallis statistical test. In the case of age information, the relative frequencies of each OTU of the total genomes and the genomes with category information were correlated using Spearman correlation. All plots were produced in R using “`ggpubr`” and `ggplot2`.

RESULTS AND DISCUSSION

Mutations Frequency Analysis

The GISAID database contains 171,461 genomes with at least 29,000 sequenced bases; from these, 109,953 genomes do not present ambiguities (as of November 30th). With an alignment of the 171,461 genomes, we performed a normalized relative frequency analysis of each nucleotide in each genomic position (NRFp; see Materials and Methods section for details). This normalization was performed to detect relevant mutations that could appear in regions where few genomes were sequenced (**Supplementary Figure S1** shows that no correlation exists between the number of cases and the number of sequenced genomes). Using this NRFp analysis, we identified nine positions estimated to be in more than 9,500,000 COVID-19 cases (more than 0.18 NRFp; **Figure 1A** and **Supplementary Figure S2A**) plus many other mutations with NRFp between 0.00 and 0.18 (**Supplementary Figures S2B,C**).

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

OTUs Identification

After NRFp analysis, we estimated a maximum likelihood tree using the whole-genome alignment of the 109,953 complete

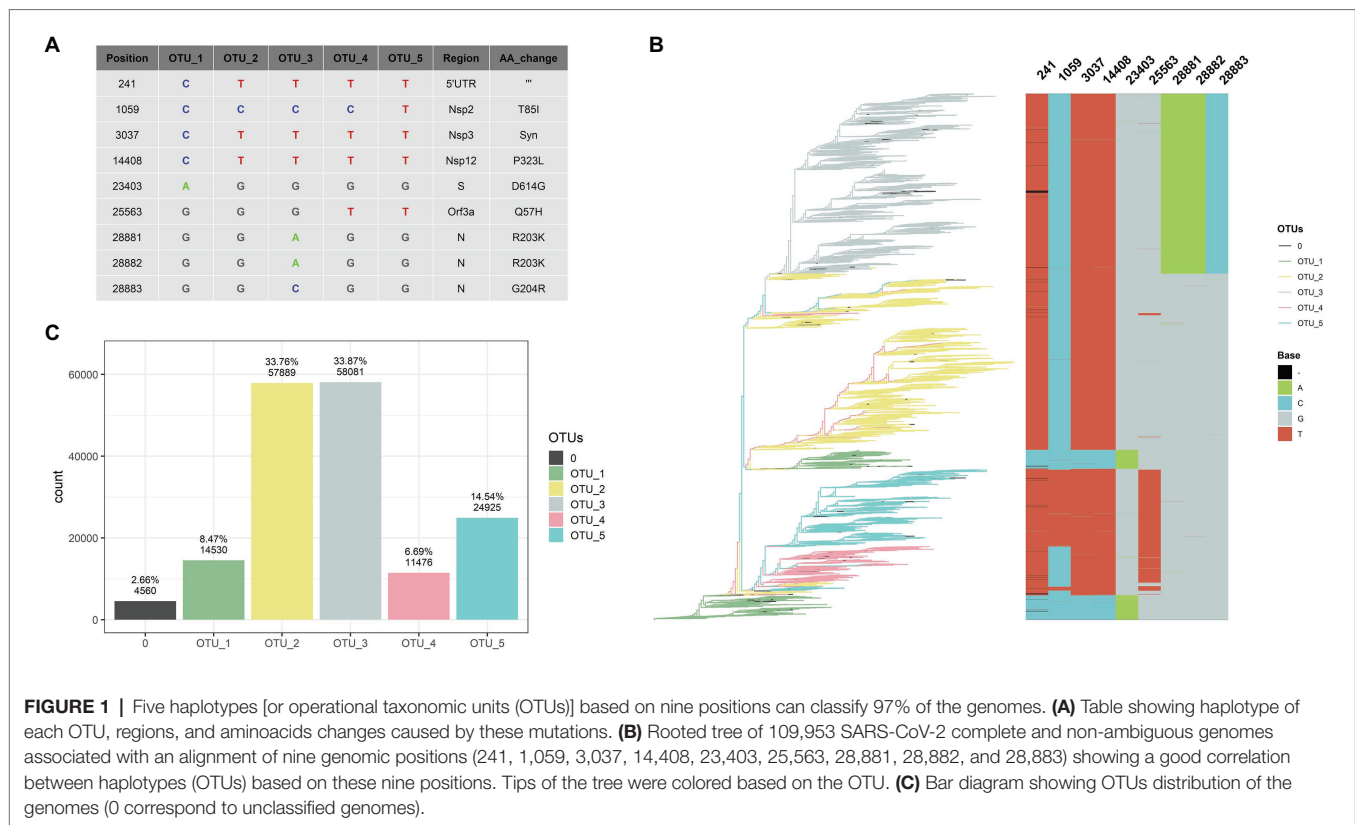


FIGURE 1 | Five haplotypes [or operational taxonomic units (OTUs)] based on nine positions can classify 97% of the genomes. **(A)** Table showing haplotype of each OTU, regions, and aminoacids changes caused by these mutations. **(B)** Rooted tree of 109,953 SARS-CoV-2 complete and non-ambiguous genomes associated with an alignment of nine genomic positions (241, 1,059, 3,037, 14,408, 23,403, 25,563, 28,881, 28,882, and 28,883) showing a good correlation between haplotypes (OTUs) based on these nine positions. Tips of the tree were colored based on the OTU. **(C)** Bar diagram showing OTUs distribution of the genomes (0 correspond to unclassified genomes).

genomes without ambiguities. Then, we associated the branches of the tree with an alignment of the nine positions (241, 1,059, 3,037, 14,408, 23,403, 25,563, 28,881, 28,882, and 28,883). We noted that combinations of those nine positions represent five well-defined groups in the tree (**Figure 1B**). Using these combinations, we defined five haplotypes that allow us to classify more than 97% of the analyzed genomes (**Figure 1C**), a great part of the remaining not classified genomes are due to the absence of sequencing corresponding to position 241. We named these haplotypes as OTUs.

OTU_1 was considered the ancestor haplotype due to its identity with the first isolated genomes (EPI_ISL_402125 and EPI_ISL_406801) with characteristic C241, C3037, C14408, and A23403. This OTU_1 comprised genomes with T or C in position 8,782 and C or T in 28,144. In other analyses, these mutations divide SARS-CoV-2 strains into two lineages. For instance, at the beginning of the pandemic, Tang et al. (2020) showed linkage disequilibrium between those positions and named them as S and L lineages. Rambaut et al. (2020) used these positions to discriminate between their proposed major lineages A and B. Those mutations did not reach the estimated number of 9,500,000 COVID-19 cases, indicating that a small number of these genomes emerged during the pandemic in comparison with other variations.

A SARS-CoV-2 isolated on January 25th in Australia is at present the first belonging to OTU_2 (**Supplementary Figure S3**). Showing simultaneously four mutations different to OTU_1 (C241T, C3037T, C14408T, and A23403G), OTU_2 is the first group containing the D614G and the P323L mutations in the

spike and nsp12 protein, respectively. Korber et al. (2020) analyzed the temporal and geographic distribution of this mutation separating SARS-CoV-2 into two groups, those with D614 and those with G614. Tomaszewski et al. (2020) analyzed the entropy of variation of these two mutations (D614G and P323L) until May. Apparently, OTU_2 is the ancestor of two other OTUs (OTU_3 and OTU_4), as shown in the maximum likelihood tree (**Figure 1B**). OTU_2 is divided into two major branches, one that originates OTU_3 and another more recent branch characteristic from Europe (see below, worldwide geographical distribution of OTUs).

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

These nine mutations have been separately described in other reports but, to our knowledge, they have not yet used been used together to classify SARS-CoV-2 haplotypes during the pandemic. The change of relative frequencies of those mutations analyzed individually showed that just in few cases, mutations that define haplotypes described here appear independently (**Supplementary Figure S4**). For example, the four mutations that define OTU_2 (C241T, C3037T, C14408T, and A23403G) rarely had been described separately and

similarly with mutations that characterize OTU_3 (G28881A, G28882A, and G28883C; **Supplementary Figure S4**). Thus, in this case, analysis of haplotypes will be identical results that if we analyzed those mutations independently.

The fact that we were able to classify more than 97% of the complete genomes data set (**Figure 1C**) shows that, at least to the present date, this classification system covers almost all the currently known genomic information around the world. Also, most of the unclassified tips appear within a clade allowing us to easily establish their phylogenetic relationships to a haplotype. Thus, at the moment, this system can be of practical use to analyze the geographical and temporal distribution of haplotypes during these 11 months of 2020. For convenience, we presented **Supplementary Table S2** that contains the relation between our identified OTUs and their relationships with pangolin lineages (Rambaut et al., 2020) and GISAID clades (Shu and McCauley, 2017).

Worldwide Geographic Distribution of OTUs

Using our OTUs classification, we analyzed the worldwide geographic distribution during 11 months of 2020. We began by plotting continental information in the ML tree of the unambiguous complete genomes (**Figure 2A**) and observed some interesting patterns. For instance, all continents contain all OTUs; also, is relatively clear that most isolates belonging to OTU_5 come from North America (**Figure 2A**). Furthermore, the biggest branch of OTU_2 is almost exclusively filled by genomes from Europe, is interesting to note that this branch also contains genomes isolated in the last months analyzed showing its relatively recent appearance (see below, the worldwide temporal distribution of OTUs). However, this approach does not allow us to evaluate continents with less sequenced genomes (**Supplementary Figure S5A**), such as South America, Oceania, and Africa. Also, it is possible that fine differences can be found in the frequency of one OTU concerning another in each continent. These differences are not observed at this level of analysis.

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

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

South America presents a greater OTU_3 frequency (**Figure 2C**) that was established in April (**Supplementary Figure S5**). This observation correlates well with studies focused in South America that detect the establishment of D614G mutation at the end of March (mutation presents in OTU_2,

OTU_3, OTU_4 and OTU_5) and a high frequency of pangolin lineage B1.1 in Chile and in general in South America that contains the same characteristic mutations that our OTU_3 (Castillo et al., 2020; Franco-Muñoz et al., 2020). Unfortunately, few genomes are reported in South America for September, October, and November (24 genomes in total in the three months), hindering a correct analysis of frequencies in these months. Similarly, OTU_3 was most prevalent in Asia, Oceania, and Africa (**Figures 2E–G**). With other OTUs with least than 0.3 NRFp (**Figures 2E–G**). Wu et al. (2020) report high incidence of mutations that define OTU_3 in Bangladesh, Oman, Russia, Australia, and Latvia. At the haplotype level, OTU_3 presents mutations in the N protein that apparently increases the fitness of this group in comparison with OTU_2 (OTU_2 does not present mutations in N; **Figure 1A**). Thus, four of the six continents analyzed present an estimation of more than 50% COVID-19 cases with a SARS-CoV-2 with the three mutations in the N protein. We, therefore, believe that it is important to more deeply study if exists positive fitness implications for these mutations.

Europe presents an interesting pattern (**Figure 2D**), it follows a similar pattern to South America, Asia, Oceania, and Africa until July (**Supplementary Figure S6**), with OTU_3 as the predominant. Then, in August, OTU_2 increased its frequency, and since September, OTU_2 is the most prevalent in Europe (**Figure 2D**). This could be caused by the appearance of mutations in the background of OTU_2 (such as those described in Justo et al., 2020b) with greater fitness than those of OTU_3 or due to other effects (i.e., founder effects) after the relaxation of lockdown policies.

Worldwide Temporal Distribution of OTUs

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

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

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

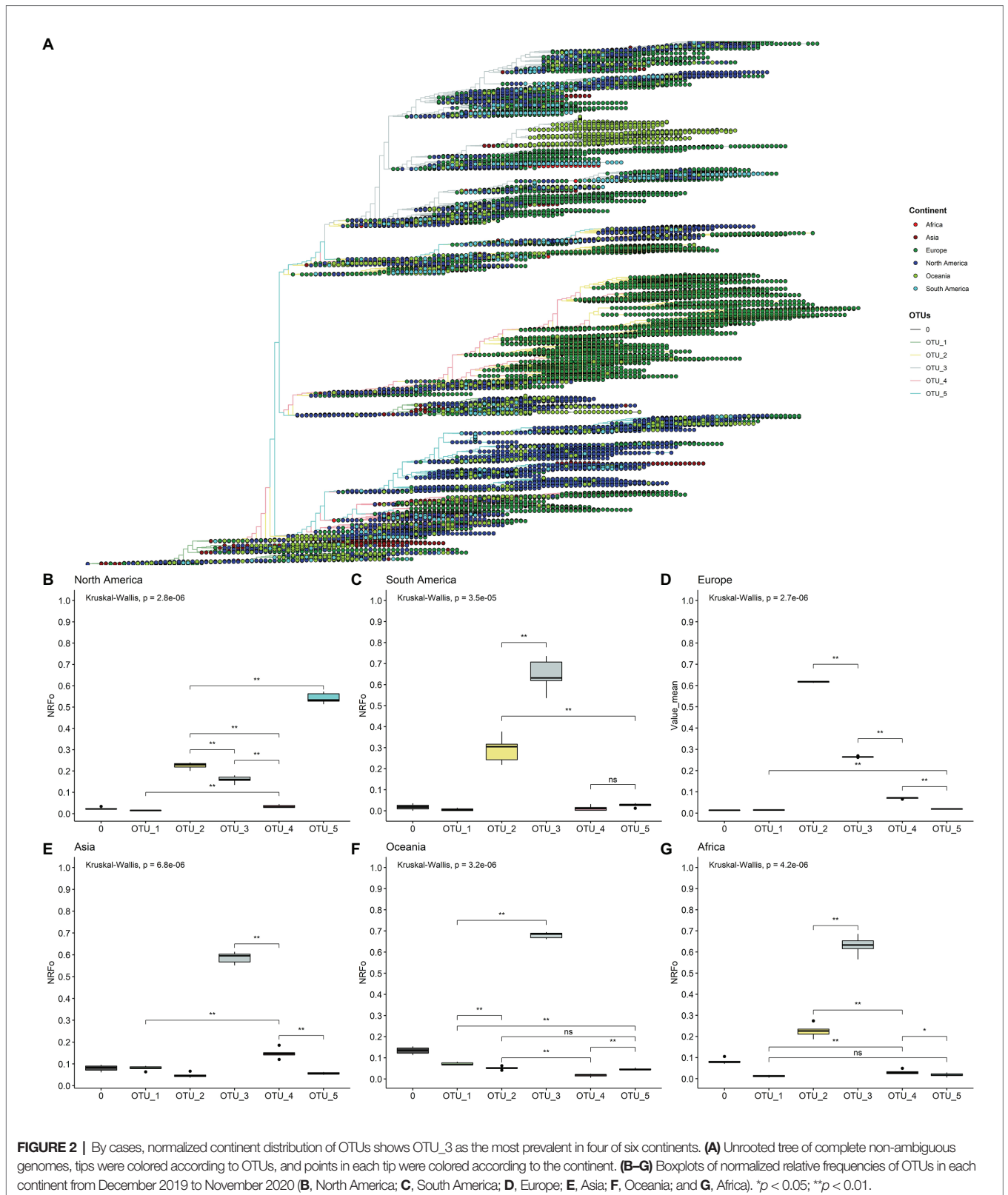
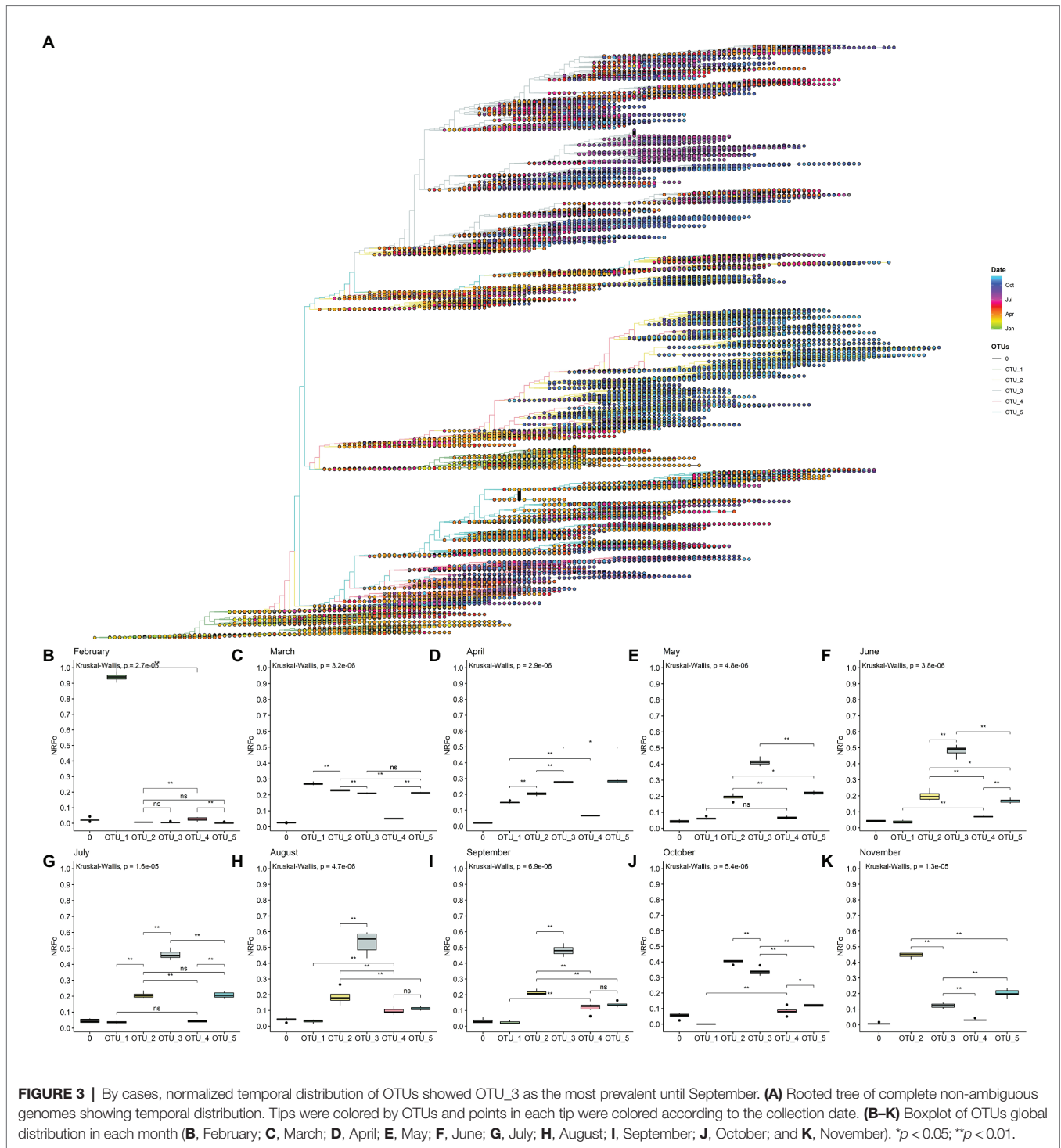


FIGURE 2 | By cases, normalized continent distribution of OTUs shows OTU_3 as the most prevalent in four of six continents. **(A)** Unrooted tree of complete non-ambiguous genomes, tips were colored according to OTUs, and points in each tip were colored according to the continent. **(B–G)** Boxplots of normalized relative frequencies of OTUs in each continent from December 2019 to November 2020 **(B, North America; C, South America; D, Europe; E, Asia; F, Oceania; and G, Africa)**. * $p < 0.05$; ** $p < 0.01$.

In March, when the epicenter of the pandemic moved to Europe and North America, but cases were still appearing in Asia, OTU_2, OTU_3, and OTU_5 increased their prevalence

but OTU_1 remained slightly as the most prevalent during this month (Figure 3C). Interestingly, OTU_4 remained in relatively low frequencies (Figure 3C). This month contains



the more homogenous OTUs distribution in a worldwide context, but with some OTUs more prevalent in each continent (Supplementary Figure S6).

During April, OTU_1 continued its downward while OTU_3 and OTU_5 increased their presence (Figure 3D) probably due to its higher representation (compared to March) in several continents such as South America, North America, and Europe (Supplementary Figure S6). During this month, Africa showed

a high prevalence of OTU_2 (Supplementary Figure S6). We also witnessed the establishment of OTU_3 in South America and OTU_5 in North America (Supplementary Figure S6).

May, June, and July showed a similar pattern, with OTU_3 as the most prevalent due to its high frequencies in South America, Oceania, and Europe (Figures 3E–G and Supplementary Figure S6). North America maintains OTU_5 as the most prevalent and Oceania showed a relatively

homogenous pattern. During these months, OTU_2 had intermediate frequencies in all continents resulting in intermediate frequencies all over the world (**Figures 3E–G** and **Supplementary Figure S6**). OTU_1 and OTU_4 representatives were reported during these months but with very low frequencies.

In August and September, we detected a slightly higher frequency of OTU_4 compared to the previous months (**Figures 3H,I**) with no significant differences with OTU_5. In September in Europe, OTU_3 stopped being the most frequent. Instead, OTU_2 was the most frequent in this month in Europe (**Supplementary Figure S6**). In October and November, OTU_2 has increased its frequency rapidly (**Figures 3J,K**) mainly due to a large number of cases and reported genomes belonging to this OTU_2 in Europe in October and November. Due to the few genomes currently available in GISAID for all continents, except for Europe and North America during November, just these two continents were analyzed in the last month.

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

Age, Gender, and Patient Status Relation With OTUs

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

Unfortunately, just 26.11% of the 171,461 genomes analyzed have age and gender information (**Supplementary Figure S8**). In the case of patient status information, we noted that GISAID categories are not well organized and we had to reclassify the information into three categories: Asymptomatic, Mild, and Severe (**Supplementary Figure S9A**). Using this classification scheme, we found that 99.14% (169,979 genomes) were not informative, 0.1% (175 genomes) falls in the Asymptomatic category, 0.33% (562 genomes) in the Mild category, and 0.43% (745 genomes) could be classified as Severe (**Supplementary Figure S9B**).

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

Thus, we concluded that none of the OTUs are related to an asymptomatic, mild, or severe COVID-19, at least in the populations analyzed.

Age information was also analyzed in the same manner. In general, although some differences were detected as significant, those were not consistently maintained between different populations analyzed (**Supplementary Figures S10A–J**). Furthermore, none difference reaches a value of p less than 0.01 (except for OTU_4 in North America). Since heterogeneity between countries information is possible, we think that these small differences are more likely due to these heterogeneities and we cannot strongly conclude that some age groups are more related to a specific OTU. Additionally, a strong positive correlation between total relative frequencies of OTUs and relative frequencies by age groups in month-continent was found, meaning that those two frequencies are similar in most of the analyzed populations (**Supplementary Figure S10K**).

A similar approach was done using gender information, but in this case, due to the greater quantity of information, we used more restrictive filter parameters. Thus, we selected country-month combinations with at least 250 genomes with male or female information and two times total genomes information (for instance USA – March has 2079 genomes from female patients and 9,287 genomes with or without gender information). Again, we did not find OTU's preference for a specific gender (**Supplementary Figure S11**).

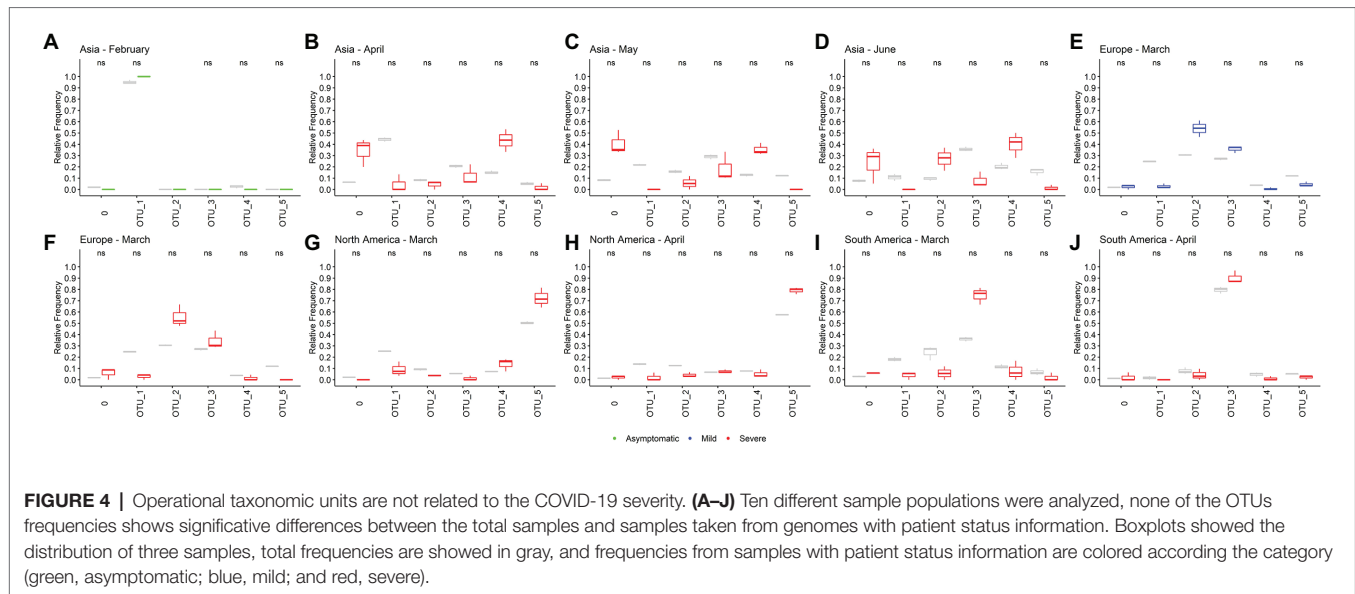
Description of the Most Frequent Mutations

C241T

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

C1059T

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



C3037T

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

C3037T causes a codon change from TTC to TTT. TTT is more frequently present in the genome of SARS-CoV-2 and other related coronaviruses compared to TTC (Gu et al., 2020) but in humans, the codon usage of TTT and TTC are similar (Mauro and Chapel, 2014). The reason why TTT is more frequent in SARS-CoV-2 is unknown but seems to be a selection related to SARS-CoV-2 and not to the host. Another option is genetic drift.

C14408T

The C14408T mutation changes P323 to leucine in Nsp12, the RNA-dependent RNA polymerase of SARS-CoV2 (Supplementary Figures S12A,B). P323 together with P322 ends helix 10 and generate a turn that is followed by a beta-sheet (Supplementary Figure S12C). Leucine at position 323 could form hydrophobic interactions with the methyl group of L324 and the aromatic ring of F396 creating a more stable variant of Nsp12 (Supplementary Figure S12E). In concordance with this, protein dynamics simulations showed a stability increase of the Nsp12 P323L variant (Chand and Azad, 2020). In the absence of P322, the mutation P323L would probably be disfavored due to the flexibilization of the turn at the end of helix 10. Experimental evidence is necessary to confirm these hypotheses and to evaluate their impact on protein function.

A23403G

An interesting protein to track is spike protein (Supplementary Figure S13A) due to its importance in SARS-CoV-2 infectivity.

It has been suggested that the D614G change in the S1 domain that results from the A23403G mutation generates a more infectious virus, less spike shedding, greater incorporation in pseudovirions (Zhang et al., 2020), and higher viral load (Korber et al., 2020).

How these effects occur at the structural level remains unclear, although some hypotheses have been put forward: (1) We think that there is no evidence for hydrogen-bond between D614 and T859 mentioned by Korber et al. (2020), and distances between D614 and T859 are too long for a hydrogen bond (Supplementary Figure S13B), (2) distances between Q613 and T859 (Supplementary Figure S13C) could be reduced by increased flexibility due to D614G substitution, forming a stabilizing hydrogen bond, and (3) currently available structures do not show salt-bridges between D614 and R646 as proposed by Zhang et al. (2020; Supplementary Figure S13D).

G25563T

Orf3a (Supplementary Figure S14A) is required for efficient *in vitro* and *in vivo* replication in SARS-CoV (Castaño-Rodríguez et al., 2018). It has been implicated in inflammasome activation (Siu et al., 2019), apoptosis (Chan et al., 2009), and necrotic cell death (Yue et al., 2018) and has been observed in Golgi membranes (Padhan et al., 2007) where pH is slightly acidic (Griffiths and Simons, 1986). Kern et al. (2020) showed that Orf3a preferentially transports Ca⁺² or K⁺ ions through a pore (Supplementary Figure S14B). Some constrictions were described in this pore, one of them formed by the side chain of Q57 (Supplementary Figure S14C).

Mutation G25563T produces the Q57H variant of Orf3a (Supplementary Figure S14C). It did not show significant differences in expression, stability, conductance, selectivity, or gating behavior (Kern et al., 2020). We modeled Q57H mutation and we did not observe differences in the radius of constriction (Supplementary Figure S14C) formed by residue 57 but

we observed slight differences in the electrostatic surface due to the ionizability of the histidine side chain (**Supplementary Figure S14D**).

G28881A, G28882A, and G28883C

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

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

CONCLUDING REMARKS

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

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

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

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

In the next months, these haplotypes description will need to be updated, identification of new haplotypes could be performed by combining the identification of new frequent mutations and phylogenetic inference. We will continue monitoring the emergence of mutations that exceed our proposed cut-off of 0.18 NRFp and this information will be rapidly shared with the scientific community through our web page.² This will also be accompanied by a continuous update of haplotypes information. During the peer-review process of this manuscript, we identify several other mutations near to the cut-off proposed that were reported in Justo et al. (2020b).

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

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

DATA AVAILABILITY STATEMENT

Publicly available datasets were analyzed in this study. This data can be found at: gisaid.org.

AUTHOR CONTRIBUTIONS

SA, DS, CH, GB, AC, and RG-SC conceived, initiated, and coordinated the project. SA performed the phylogenetic analyses, geographical and temporal analyses. GU-C wrote python scripts used in data processing and analyses. SA, DS, CH, GB, AC, RG-SC and RC performed the structural analysis. The manuscript was written by SA, DS, CH, and GB. All authors discussed the methodologies and results, and read and approved the manuscript.

FUNDING

We thank to the Fundação de Amparo à Pesquisa do Estado de São Paulo (FAPESP) graduate scholarship (to SA; 2015/13318-4) and Universidad Ricardo Palma (URP) for APC financing.

²<https://sarscov2haplofinder.urp.edu.pe/>

ACKNOWLEDGMENTS

This manuscript has been released as a pre-print at <https://doi.org/10.1101/2020.07.12.199414> (Justo et al., 2020a).

We are very grateful to the GISAID Initiative and all its data contributors, i.e., the authors from the Originating laboratories responsible for obtaining the specimens and the Submitting laboratories where genetic sequence data were generated and shared *via* the GISAID Initiative, on which this research is based. Complete acknowledgments of the 171,461 genomes used are available in **Supplementary Material (SF1–SF20)**.

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

Pinho (Albert Einstein Hospital – Sao Paulo) and Deyvid Amgarten (Albert Einstein Hospital – Sao Paulo) for its helpful comments. To the Ricardo Palma University High-Performance Computational Cluster (URPHPC) managers Gustavo Adolfo Abarca Valdiviezo and Roxana Paola Mier Hermoza at the Ricardo Palma Informatic Department (OFICIC) for their contribution in programs and remote use configuration of URPHPC. To Gladys Arevalo Chong for her figure style suggestions.

SUPPLEMENTARY MATERIAL

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fmicb.2021.612432/full#supplementary-material>

REFERENCES

- Blom, N., Sicheritz-Pontén, T., Gupta, R., Gammeltoft, S., and Brunak, S. (2004). Prediction of post-translational glycosylation and phosphorylation of proteins from the aminoacid sequence. *Proteomics* 4, 1633–1649. doi: 10.1002/pmic.200300771
- Castañero-Rodríguez, C., Honrubia, J., Gutierrez-Alvarez, J., DeDiego, M., Nieto-Torres, J., Jimenez-Guardeño, J., et al. (2018). Role of severe acute respiratory syndrome coronavirus viroporins E, 3a, and 8a in replication and pathogenesis. *mBio* 9, e02325–e02417. doi: 10.1128/mBio.02325-17
- Castillo, A. E., Parra, B., Tapia, P., Lagos, J., Arata, L., Acevedo, A., et al. (2020). Geographical distribution of genetic variants and lineages of SARS-CoV-2 in Chile. *Front. Public Health* 8:562615. doi: 10.3389/fpubh.2020.562615
- Chan, C., Tsoi, H., Chan, W., Zhai, S., Wong, C., Yao, X., et al. (2009). The ion channel activity of the SARS-coronavirus 3a protein is linked to its pro-apoptotic function. *Int. J. Biochem. Cell Biol.* 41, 2232–2239. doi: 10.1016/j.biocel.2009.04.019
- Chand, G., and Azad, G. (2020). Identification of novel mutations in RNA-dependent RNA polymerases of SARS-CoV-2 and their implications. *bioRxiv* [Preprint]. doi: 10.1101/2020.05.05.079939
- Chang, C., Chen, C., Chiang, M., Hsu, Y., and Huang, T. (2013). Transient oligomerization of the SARS-CoV N protein – Implication for virus ribonucleoprotein packaging. *PLoS One* 8:e65045. doi: 10.1371/journal.pone.0065045
- Chang, C., Hsu, Y., Chang, Y., Chao, F., Wu, M., Huang, Y., et al. (2009). Multiple nucleic acid binding sites and intrinsic disorder of severe acute respiratory syndrome coronavirus nucleocapsid protein implications for ribonucleocapsid protein packaging. *J. Virol.* 83, 2255–2264. doi: 10.1128/JVI.02001-08
- Cornillez-Ty, C., Liao, L., Yates, J., Kuhn, P., and Buchmeier, M. (2009). Severe acute respiratory syndrome coronavirus nonstructural protein 2 interacts with a host protein complex involved in mitochondrial biogenesis and intracellular signaling. *J. Virol.* 83, 10314–10318. doi: 10.1128/JVI.00842-09
- Cucinotta, D., and Vanelli, M. (2020). WHO declares COVID-19 a pandemic. *Acta Biomed.* 91, 157–160. doi: 10.23750/abm.v91i1.9397
- Franco-Muñoz, C., Álvarez-Díaz, D., Laiton-Donato, K., Wiesner, M., Escandón, P., Usme-Ciro, J., et al. (2020). Substitutions in spike and nucleocapsid proteins of SARS-CoV-2 circulating in South America. *Infect. Genet. Evol.* 85:104557. doi: 10.1016/j.meegid.2020.104557
- Fusaro, G., Dasgupta, P., Rastogi, S., Joshi, B., and Chellappan, S. (2003). Prohibitin induces the transcriptional activity of p53 and is exported from the nucleus upon apoptotic signaling. *J. Biol. Chem.* 278, 47853–47861. doi: 10.1074/jbc.M305171200
- Graham, R., Sims, A., Brockway, S., Baric, S., and Denison, M. (2005). The nsp2 replicase protein of murine hepatitis virus and severe acute respiratory syndrome coronavirus is dispensable for viral replication. *J. Virol.* 79, 13399–13411. doi: 10.1128/JVI.79.21.13399-13411.2005
- Pinho (Albert Einstein Hospital – Sao Paulo) and Deyvid Amgarten (Albert Einstein Hospital – Sao Paulo) for its helpful comments. To the Ricardo Palma University High-Performance Computational Cluster (URPHPC) managers Gustavo Adolfo Abarca Valdiviezo and Roxana Paola Mier Hermoza at the Ricardo Palma Informatic Department (OFICIC) for their contribution in programs and remote use configuration of URPHPC. To Gladys Arevalo Chong for her figure style suggestions.
- Griffiths, G., and Simons, K. (1986). The trans Golgi network: sorting at the exit site of the golgi complex. *Science* 234, 438–443. doi: 10.1126/science.2945253
- Gu, H., Chu, D., Peiris, M., and Poon, L. (2020). Multivariate analyses of codon usage of SARS-CoV-2 and other betacoronaviruses. *bioRxiv* [Preprint]. doi: 10.1101/2020.02.15.950568
- Justo, S., Zapata, D., Huallpa, C., Landa, G., Castillo, A., Garavito-Salini, R., et al. (2020a). Global geographic and temporal analysis of SARS-CoV-2 haplotypes normalized by COVID-19 cases during the pandemic. *bioRxiv* [Preprint]. doi: 10.1101/2020.07.12.199414
- Justo, S., Zapata, D., Huallpa, C., Landa, G., Castillo, A., Garavito-Salini, R., et al. (2020b). Analysis of the dynamics and distribution of SARS-CoV-2 mutations and its possible structural and functional implications. *bioRxiv* [Preprint]. doi: 10.1101/2020.11.13.381228
- Kepler, L., Hamins-Puertolas, M., and Rasmussen, D. (2020). Decomposing the sources of SARS-CoV-2 fitness variation in the United States. *bioRxiv* [Preprint]. doi: 10.1101/2020.12.14.422739
- Kern, D., Sorum, B., Hoel, C., Sridharan, S., Remis, J., Toso, D., et al. (2020). Cryo-EM structure of the SARS-CoV-2 3a ion channel in lipid nanodiscs. *bioRxiv* [Preprint]. doi: 10.1101/2020.06.17.156554
- Korber, B., Fischer, W., Gnanakaran, S., Yoon, H., Theiler, J., Abfalterer, W., et al. (2020). Tracking changes in SARS-CoV-2 Spike: evidence that D614G increases infectivity of the COVID-19 virus. *Cell* 182, 812.e19–827.e19. doi: 10.1016/j.cell.2020.06.043
- Kuo, P., Chiang, C., Wang, Y., Doudeva, L., and Yuan, H. (2014). The crystal structure of TDP-43 RRM1-DNA complex reveals the specific recognition for UG- and TG-rich nucleic acids. *Nucleic Acids Res.* 42, 4712–4722. doi: 10.1093/nar/gkt1407
- Lee, E., Lee, V., and Trojanowski, J. (2011). Gains or losses: molecular mechanisms of TDP43-mediated neurodegeneration. *Nat. Rev. Neurosci.* 13, 38–50. doi: 10.1038/nrn3121
- Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 34, 3094–3100. doi: 10.1093/bioinformatics/bty191
- Luo, H., Chen, Q., Chen, J., Chen, K., Shen, X., and Jiang, H. (2005). The nucleocapsid protein of SARS coronavirus has a high binding affinity to the human cellular heterogeneous nuclear ribonucleoprotein A1. *FEBS Lett.* 579, 2623–2628. doi: 10.1016/j.febslet.2005.03.080
- Madhugiri, R., Fricke, M., Marz, M., and Ziebuhr, J. (2014). RNA structure analysis of alphacoronavirus terminal genome regions. *Virus Res.* 194, 76–89. doi: 10.1016/j.virusres.2014.10.001
- Maitra, A., Chawla, M., Raheja, H., Biswas, N., Chakraborti, S., Kumar, A. M., et al. (2020). Mutations in SARS-CoV-2 viral RNA identified in Eastern India: possible implication for the ongoing outbreak in India and impact on viral structure and host susceptibility. *J. Biosci.* 45:76. doi: 10.1007/s12038-020-00046-1
- Masters, P. (2019). Coronavirus genomic RNA packaging. *Virology* 537, 198–207. doi: 10.1016/j.virol.2019.08.031

- Mauro, V., and Chapel, S. (2014). A critical analysis of codon optimization in human therapeutics. *Trends Mol. Med.* 20, 604–613. doi: 10.1016/j.molmed.2014.09.003
- Mercatelli, D., and Giorgi, F. (2020). Geographic and genomic distribution of SARS-CoV-2 mutations. *Front. Microbiol.* 11:1800. doi: 10.3389/fmicb.2020.01800
- Merkwirth, C., and Langer, T. (2008). Prohibitin function within mitochondria: essential roles for cell proliferation and cristae morphogenesis. *Biochim. Biophys. Acta* 1793, 27–32. doi: 10.1016/j.bbamcr.2008.05.013
- Moshiri, N. (2020). ViralMSA: massively scalable reference-guided multiple sequence alignment of viral genomes. *Bioinformatics* btaa743. doi: 10.1093/bioinformatics/btaa743 [Epub ahead of print]
- Mukherjee, M., and Goswami, S. (2020). Global cataloging of variations in untranslated regions of viral genome and prediction of key host RNA binding protein-microRNA interactions modulating genome stability in SARS-CoV-2. bioRxiv [Preprint]. doi: 10.1101/2020.06.09.134585
- Pachetti, M., Marini, B., Benedetti, F., Giudici, F., Mauro, E., Storici, P., et al. (2020). Emerging SARS-CoV-2 mutation hot spots include a novel RNA-dependent-RNA polymerase variant. *J. Transl. Med.* 18:179. doi: 10.1186/s12967-020-02344-6
- Padhan, K., Tanwar, C., Hussain, A., Hui, P., Lee, M., Cheung, C., et al. (2007). Severe acute respiratory syndrome coronavirus Orf3a protein interacts with caveolin. *J. Gen. Virol.* 88, 3067–3077. doi: 10.1099/vir.0.82856-0
- Price, M., Dehal, P., and Arkin, A. (2009). FastTree: computing large minimum-evolution trees with profiles instead of a distance matrix. *Mol. Biol. Evol.* 26, 1641–1650. doi: 10.1093/molbev/msp077
- Price, M., Dehal, P., and Arkin, A. (2010). FastTree 2 – approximately maximum-likelihood trees for large alignments. *PLoS One* 5:e9490. doi: 10.1371/journal.pone.0009490
- Rahman, M., Islam, M., Alam, A., Islam, I., Hoque, M., Akter, S., et al. (2020). Evolutionary dynamics of SARS-CoV-2 nucleocapsid protein and its consequences. *J. Med. Virol.* 1–19. doi: 10.1002/jmv.26626 [Epub ahead of print]
- Rajalingam, K., Wunder, C., Brinkmann, V., Churin, Y., Hekman, M., Sievers, C., et al. (2005). Prohibitin is required for RAS-induced RAF-MEK-ERK activation and epithelial cell migration. *Nat. Cell Biol.* 7, 837–843. doi: 10.1038/ncb1283
- Rambaut, A., Holmes, E., Hill, V., O’Toole, A., Hill, V., McCrone, J., et al. (2020). A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology. *Nat. Microbiol.* 5, 1403–1407. doi: 10.1038/s41564-020-0770-5
- Saha, O., Hossain, M., and Rahaman, M. (2020). Genomic exploration light on multiple origin with potential parsimony-informative sites of the severe acute respiratory syndrome coronavirus 2 in Bangladesh. *Gene Rep.* 21:100951. doi: 10.1016/j.genrep.2020.100951
- Shu, Y., and McCauley, J. (2017). GISAID: global initiative on sharing all influenza data – from vision to reality. *Euro Surveill.* 22, 1–3. doi: 10.2807/1560-7917.ES.2017.22.13.30494
- Siu, K., Yuen, K., Castaño-Rodríguez, C., Ye, Z., Yeung, M., Fung, S., et al. (2019). Severe acute respiratory syndrome coronavirus ORF3a protein activates the NLRP3 inflammasome by promoting TRAF3-dependent ubiquitination of ASC. *FASEB J.* 33, 8865–8877. doi: 10.1096/fj.201802418R
- Sun, L., Liu, L., Yang, X., and Wu, Z. (2004). Akt binds prohibitin 2 and relieves its repression of MyoD and muscle differentiation. *J. Cell Sci.* 117, 3021–3029. doi: 10.1242/jcs.01142
- Tang, X., Wi, C., Li, X., Song, Y., Yao, X., Wu, X., et al. (2020). On the origin and continuing evolution of SARS-CoV-2. *Natl. Sci. Rev.* 7, 1012–1023. doi: 10.1093/nsr/nwaa036
- Tomaszewski, T., DeVries, R., Dong, M., Bhatia, G., Norsworthy, M., Zheng, X., et al. (2020). New pathways of mutational change in SARS-CoV-2 proteomes involve regions of intrinsic disorder important of virus replication and release. *Evol. Bioinform.* 16, 1–18. doi: 10.1177/1176934320965149
- van Dorp, L., Richard, D., Tan, C., Shaw, L., Acman, M., and Balloux, F. (2020). No evidence for increased transmissibility from recurrent mutations in SARS-CoV-2. *Nat. Commun.* 11:5986. doi: 10.1038/s41467-020-19818-2
- Wang, S., Nath, N., Adlam, M., and Chellappan, S. (1999). Prohibitin, a potential tumor suppressor, interacts with RB and regulates E2F function. *Oncogene* 18, 3501–3510. doi: 10.1038/sj.onc.1202684
- World Health Organization (2020). Available at: <https://covid19.who.int/> (Accessed August 25, 2020).
- Wu, S., Tian, C., Liu, P., Guo, D., Zheng, W., Huang, X., et al. (2020). Effects of SARS-CoV-2 mutations on protein structures and intraviral protein-protein interactions. *J. Med. Virol.* doi: 10.1002/jmv.26597 [Epub ahead of print]
- Yu, G. (2020). Using ggtree to visualize data on tree-like structures. *Curr. Protoc. Bioinformatics* 69, 1–18. doi: 10.1002/cpbi.96
- Yu, G., Smith, D., Zhu, H., Guan, Y., and Lam, T. (2017). GGTREE: an R package for visualization and annotation of phylogenetic trees with their covariates and other associated data. *Methods Ecol. Evol.* 8, 28–36. doi: 10.1111/2041-210X.12628
- Yue, Y., Nabar, N., Shi, C., Kamenyeva, O., Xiao, X., Hwang, I., et al. (2018). SARS-Coronavirus open reading frame-3a drives multimodal necrotic cell death. *Cell Death Dis.* 9, 1–15. doi: 10.1038/s41419-018-0917-y
- Yun, C. (2020). Genotyping coronavirus SARS-CoV-2: methods and implication. *Genomics* 112, 3588–3596. doi: 10.1016/j.ygeno.2020.04.016
- Zhang, L., Jackson, C., Mou, H., Ojha, A., Rangarajan, E., Izard, T., et al. (2020). The D614G mutation in the SARS-CoV-2 spike protein reduces S1 shedding and increases infectivity. bioRxiv [Preprint]. doi: 10.1101/2020.06.12.148726
- Zhu, N., Zhang, D., Wang, W., Li, X., Yang, B., Song, J., et al. (2020). A novel coronavirus from patients with pneumonia in China, 2019. *N. Engl. J. Med.* 382, 727–733. doi: 10.1056/NEJMoa2001017

Conflict of Interest: The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2021 Justo Arevalo, Zapata Sifuentes, Huallpa, Landa Bianchi, Castillo Chávez, Garavito-Salini Casas, Uceda-Campos and Pineda Chavarria. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.