

# 1 **Positive selection of ORF3a and ORF8 genes drives the evolution of** 2 **SARS-CoV-2 during the 2020 COVID-19 pandemic**

3 **Lauro Velazquez-Salinas**<sup>1,2 \*</sup>, **Selene Zarate**<sup>3</sup>, **Samantha Eberl**<sup>4</sup>, **Douglas P**<sup>1</sup>. **Gladue,**  
4 **Isabel Novella**<sup>5</sup> and **Manuel V. Borca**<sup>1</sup>

5 Affiliations

6 1. Foreign Animal Disease Research Unit, USDA/ARS Plum Island Animal Disease Center, PO  
7 Box 848, Greenport NY 11944, USA.

8 2. College of Veterinary Medicine, Kansas State University, Manhattan, KS 66506, USA.

9  
10 3. Posgrado en Ciencias Genómicas, Universidad Autónoma de la Ciudad de Mexico, Mexico  
11 City, Mexico.

12 4. Central Connecticut State University, New Britain, CT 06053.

13 5. Independent researcher

14 \* Corresponding author: Lauro Velazquez-Salinas. [Lauro.velazquez@usda.gov](mailto:Lauro.velazquez@usda.gov)

15

## 16 **Abstract**

17 In this study, we analyzed full-length SARS-CoV-2 genomes from multiple countries to  
18 determine early trends in the evolutionary dynamics of the novel COVID-19 pandemic. Results  
19 indicated SARS-CoV-2 evolved early into at least three phylogenetic groups, characterized by  
20 positive selection at specific residues of the accessory proteins OFR3a and ORF8a. We also  
21 report evidence of epistatic interactions among sites in the genome that may be important in the  
22 generation of variants adapted to humans. These observations might impact not only public  
23 health, but also suggest more studies are needed to understand the genetic mechanisms that may  
24 affect the development of therapeutic and preventive tools, like antivirals and vaccines.

25 **Keywords:** Novel coronavirus, pandemic, evolution, epistasis, positive selection, COVID-19,  
26 SARS-CoV2, phylogenetics.

27

## 28 **Introduction.**

29 The first case of pneumonia confirmed to be caused by the novel virus SARS-CoV-2 was in a  
30 41-year-old man in Wuhan, Hubei province, China on 31/December/2019 (Wu et al., 2020). As  
31 of 9/April/2020, the World Health Organization (WHO) has confirmed 1,439,516 cases, 85,711  
32 deaths, and the presence of COVID-19 in 209 countries, areas or territories. Of the confirmed  
33 cases, 71% are from seven countries: United States of America (395,030), Spain (146, 690),  
34 Italy (139, 422), Germany (108,202), China (83,249), France (81,095), and Iran (66,220)  
35 (<https://www.who.int/emergencies/diseases/novel-coronavirus-2019>). As of the writing of this  
36 report, the number of COVID-19 cases continue to increase worldwide, with multiple epicenters.

37 The International Committee on Taxonomy of Viruses (ICTV) initially named this pathogen  
38 2019-nCoV (also referred to as COVID-19 by WHO) and included it within the *Coronaviridae*  
39 viral family (Coronaviridae Study Group of the International Committee on Taxonomy of, 2020).  
40 Later, based on the close phylogenetic relationship of COVID-19 with other human and bat  
41 severe acute respiratory syndrome coronaviruses (SARS-CoVs), ICTV renamed the virus as  
42 SARS-CoV-2 (Coronaviridae Study Group of the International Committee on Taxonomy of,  
43 2020).

44 The *Coronaviridae* family encompasses a group of single-stranded, positive-sense RNA viruses  
45 with a genome length varying between 27 and 32 kb. These are zoonotic viruses with the  
46 potential to infect humans and animals. Coronaviruses may cause acute and chronic respiratory,  
47 enteric, and central nervous system infections (Phan et al., 2018, Weiss and Navas-Martin, 2005).  
48 In the case of SARS-CoV-2, a meta-analysis of 50,466 patients indicate that fever and cough are  
49 the most common symptoms (95% CI: 81.8-94.5% and 65.7-78.2%, respectively) (Sun et al.,  
50 2020). The disease may worsen, and the percentages of severe cases and fatality rate vary  
51 between 12.7-24.3% and 2.7-6.1% (95% CI), respectively (Sun et al., 2020).

52 The genome organization of SARS-CoV-2 is similar to viruses from the genus *Betacoronavirus*,  
53 one of the four genera included in the *Coronaviridae* subfamily Orthocoronavirinae. The  
54 ~29,903 nucleotide (nt) genome is organized as follows, 5' to 3': replicase ORF1ab, S (encoding  
55 the structural spike glycoprotein), ORF3a (ORF3a protein), E (structural envelope protein), M  
56 (structural membrane glycoprotein), ORF6 (ORF6 protein), ORF7a (ORF7a protein), ORF7b  
57 (ORF7b protein), ORF8 (ORF8 protein), N (structural nucleocapsid phosphoprotein), and  
58 ORF10 (ORF10 protein). ORF1ab (~21,291 nt) encodes sixteen non-structural proteins: leader  
59 protein, nsp2, nsp3, nsp4, 3C-like proteinase, nsp6, nsp7, nsp8, nsp9, nsp10, RNA-dependent  
60 RNA polymerase, helicase, 3'-to 5' exonuclease, endoRNase, 2'-o-ribose methyltransferase, and  
61 nsp11 (Wu et al., 2020).

62 Much speculation regarding the origin of SARS-CoV-2 emanates from unfounded theories, such  
63 as a man-made laboratory origin; however, a recent study supports the hypothesis that SARS-  
64 CoV-2 was the result of cross-species transmission followed by natural selection in the novel  
65 human host (Andersen et al., 2020). This hypothesis is strongly supported by studies examining  
66 amino acid differences between SARS-CoV-2 and some phylogenetically related  
67 betacoronaviruses (e.g., Bat-RatG13 isolate and the human SARS-CoV isolate Urbani) at the  
68 receptor-binding domain (RBD) of the spike protein, where such differences seem to increase  
69 the ability of SARS-CoV-2 to bind to the human receptor angiotensin-converting enzyme 2  
70 (ACE2) (Andersen et al., 2020). This increased affinity for binding ACE2 might help to explain  
71 the infectiousness of SARS-CoV-2 in human populations (Wan et al., 2020).

72 Considering the extraordinary plasticity shown by other human viral RNA pathogens, for  
73 example HIV-1, Influenza viruses, SARS-CoV, and hepatitis C virus, to undergo adaptative  
74 changes to evade innate and adaptive immune responses, develop drug resistance or establish an  
75 infection in a new host (Frost et al., 2018), multiple questions arise regarding the adaptative  
76 changes that SARS-CoV-2 has undergone during the pandemic. SARS-CoV-2 has spread  
77 throughout many countries resulting in the infection of people with diverse immunological  
78 backgrounds and demographics (age, sex, environmental conditions, etc) that potentially impose  
79 significant selective pressures on SARS-CoV-2.

80 Here, we evaluate the phylogenetic and evolutionary dynamics of SARS-CoV-2 during the early  
81 phase of the COVID-19 pandemic. Using different analyses based on a codon-based  
82 phylogenetic framework, we identified critical sites in the genome undergoing positive selection,  
83 which might favor viral divergence and emergence of multiple viral variants. Our findings are  
84 discussed in terms of the potential effects that the rapid evolution of SARS-CoV-2 might have on  
85 public health.

## 86 **Materials and methods**

### 87 **Data collection**

88 Eighty-six full-length SARS-CoV-2 genomes representing early viral isolates from patients  
89 living in diverse geographic regions were used for this study. Viral sequences, downloaded from  
90 the NCBI GenBank database on 06/03/2020, represent the total number of full-length viral  
91 genomes at the time that the analysis was conducted (Figure 1).

### 92 **Phylogenetic dynamic analysis**

93 Two phylogeographic trees were reconstructed using the programs BEAST2.4.3, BEAUti, and  
94 TreeAnnotator (Drummond et al., 2012), introducing the country/region as a trait. The Markov  
95 Chain Monte Carlo was run for 25 million generations, using the HKY85 substitution model and  
96 a gamma distribution with four categories as the site heterogeneity model. The resulting file was  
97 analyzed with Tracer 1.6 to check for convergence and to determine the Burnin proportion  
98 (Rambaut et al., 2018). Finally, TreeAnnotator was used to build the maximum clade credibility  
99 tree, which was visualized with FigTree 1.4.3  
100 (<http://tree.bio.ed.ac.uk/software/figtree>). Additionally, the evolutionary rate of SARS-CoV-2,  
101 expressed as substitutions/sites/year, was calculated using the same methodology but using  
102 sampling date as a trait.

### 103 **Pairwise distance calculations**

104 Nucleotide and amino acid pairwise distance calculations among SARS-CoV-2 sequences were  
105 conducted using the SSE 1.3 Sequence Distances program (Simmonds, 2012). For this purpose, a  
106 sliding window of 50 nucleotides (nt), with a shift of 25 nt, was used to determine pairwise  
107 distances. Additionally, p-distances in nucleotide and amino acid sequences between  
108 phylogenetic groups were calculated using MEGA 7 (Kumar et al., 2016).

### 109 **Evolutionary rate per site analysis**

110 Mean (relative) evolutionary rates for each site in the alignment were estimated under the  
111 General Time Reversible model, including all 3 codon positions. These rates were scaled  
112 considering the average evolutionary rate across all sites is 1. This means that sites showing a  
113 rate <1 are evolving slower than average, and those with a rate >1 are evolving faster than  
114 average. This analysis was conducted using MEGA 7 (Kumar et al., 2016).

### 115 **Inference of selective pressures**

116 Since natural selection can be manifested as different modes (diversifying, directional or  
117 purifying), we used a combination of different evolutionary analyses to enhance the detection of  
118 relevant sites in the genome of SARS-CoV-2 experiencing diversifying (positive) and purifying  
119 (negative) selection: Single Likelihood Ancestor Counting (SLAC) (Kosakovsky Pond and Frost,

120 2005), Fixed Effects Likelihood (FEL) (Kosakovsky Pond and Frost, 2005), Mixed Effects  
121 Model of Evolution (MEME) (Murrell et al., 2012), and Fast Unbiased Bayesian Approximation  
122 (FUBAR) (Murrell et al., 2013). These methods use a maximum likelihood or Bayesian approach  
123 (FUBAR) to infer nonsynonymous (dN) and synonymous (dS) substitution rates on a per site  
124 basis for given coding alignment and corresponding phylogeny (Weaver et al., 2018). SLAC,  
125 FEL and FUBAR were methods used to identify sites experiencing pervasive diversifying or  
126 purifying selection, while MEME was used to detect sites experiencing episodic diversifying  
127 selection.

128 The presence of recombination in the sequence dataset potentially affecting the detection of  
129 positive selection was assessed using the algorithm GARD (Kosakovsky Pond et al., 2006). All  
130 methods were performed on the adaptive evolution server Datamonkey 2.0 (Weaver et al., 2018).

131 Evidence of directional selection was assessed on amino acid sequences using the Directional  
132 Evolution of Protein Sequences (DEPS) method, implemented on the Datamonkey webserver  
133 (classic) (Delport et al., 2010). This method is a model-based phylogenetic maximum likelihood  
134 test that looks for evidence of preferential substitution toward a given residue at individual  
135 positions of a protein alignment (Kosakovsky Pond et al., 2008).

### 136 **Coevolution analysis**

137 Evidence of coevolution among different sites in the SARS-CoV-2 genome was evaluated using  
138 the method Bayesian Graphical Models for co-evolving sites (BGM) (Poon et al., 2007). This  
139 method detects coevolutionary interactions between amino acids in a protein, where amino acid  
140 substitutions are mapped to branches in the phylogenetic tree.

### 141 **Blosum 62 substitution matrix (BSM62)**

142 BSM62 was used to infer the biological significance of amino acid replacements found during  
143 the evolutionary analysis of SARS-CoV-2, where positive values reflect that the substitution is  
144 most likely a product of random substitution, while negative values are indicative of selection  
145 (Henikoff and Henikoff, 1992).

## 146 **Results**

### 147 **Phylogenetic dynamics of SARS-CoV-2**

148 To evaluate potential divergence events of SARS-CoV-2, indicating the rise of new variants  
149 early during the pandemic, we reconstructed the evolution of SARS-CoV-2 using full-length  
150 genome sequences of viruses collected between late December of 2019 and early March of 2020  
151 from patients infected in different countries around the world. The results of the phylogenetic  
152 analysis demonstrate the rapid divergence of SARS-CoV-2 into three distinct phylogenetic  
153 groups, differentiated only by a few changes at the nucleotide and amino acid levels (Figure 2A  
154 and 2B).

155 Group A includes the deduced ancestral sequence (NC\_045512.2) obtained from the index case  
156 in Wuhan, China, and reported to the WHO on 31/December/ 2019, as well as multiple viral  
157 isolates from different Chinese provinces. The position of these sequences among multiple  
158 branches within the Group A cluster suggests the emergence of multiple viral variants in China,  
159 especially from Wuhan before the start of the global pandemic. Furthermore, the basal branch  
160 position of some of these variants indicates that they were the ancestors of viral isolates obtained

161 from patients in the USA, Japan, Finland, Taiwan, Nepal, and India between January and  
162 February 2020.

163 Similarly, in the Group B cluster, we found viral isolates from multiple Chinese provinces  
164 between December of 2019 and January of 2020. These isolates are likely ancestors of viral  
165 isolates recovered from patients in the USA, India, and Taiwan between January and March  
166 2020. Interestingly, one isolate from Wuhan (LR757995.1) is part of the Group B cluster,  
167 supporting the hypothesis that multiple viral variants emerged in China before the start of the  
168 pandemic.

169 The Group C cluster was the only cluster that did not contain sequences from China. This cluster  
170 includes viral isolates collected from the USA, Italy, Australia, Sweden, Brazil, and South Korea  
171 between January and February of 2020. The absence of viral isolates from China and the  
172 increased genetic distance from Group A suggests that the emergence of these variants might  
173 have come from a second wave of transmission outside of China after the start of the pandemic.

### 174 **Evolutionary divergence in the genome of SARS-CoV-2**

175 Once we reconstructed the phylodynamic of SARS-CoV-2 isolates obtained early during the  
176 pandemic event, we attempted to determine which nucleotide positions in the SARS-CoV-2  
177 genome were related to the early divergence of this virus. Overall, the evolutionary rate of  
178 SARS-CoV-2 is  $1.15 \times 10^{-3}$  substitutions/site/year (95% HPD  $7.41 \times 10^{-4}$  -  $1.57 \times 10^{-3}$ ), while  
179 pairwise analysis at nucleotide and amino acid levels revealed an average identity of 99.93-  
180 99.98% and 99.86-99.97%, respectively. Given the short divergence time, a high level of identity  
181 is to be expected; however, a few synonymous and non-synonymous substitutions were observed  
182 in the ORF1ab, S, ORF3a, M, ORF8a, N and ORF10 genes (Figure 3A and 3B). When pairwise  
183 distances were calculated based on gene length, the highest levels of divergence were observed  
184 within genes ORF10 and ORF8a when considering synonymous and non-synonymous  
185 substitutions, respectively (Figure 3C and 3D).

186 Also, the estimated per site evolutionary rate in the coding regions revealed that 98.85% of the  
187 sites in the genome are evolving at expected rates of evolution, while 1.15% of the sites are  
188 evolving faster than expected (Figure 3E). In this context, and consistent with the length of the  
189 ORF1ab gene, most of these synonymous and non-synonymous substitutions (82 sites) were  
190 distributed among different protein-encoding segments of this gene; the segment encoding nsp3  
191 had the highest number of polymorphic sites (Figure 3E).

### 192 **Detection of purifying and diversifying selection**

193 Once we identified fast-evolving positions within different genes of SARS-CoV-2, we used a  
194 combination of different algorithms centered on a codon-based phylogenetic framework to detect  
195 specific codons evolving under natural selection. Overall, no recombination events potentially  
196 affecting the results of these analyses were detected using the GARD algorithm.

197 Using SLAC we obtained a broad picture of the extent of natural selection acting upon the  
198 SARS-CoV-2 genome. We found an overall dN/dS ratio of 0.937 along the genome. In  
199 particular, 75 codons located within 5 genes (ORF1ab > S > N > ORF8a > ORF3a) showed  
200 evidence of increased fixation of non-synonymous mutations (dN/dS > 1). Conversely, a small  
201 number of codons (35 codons) located within 5 genes (ORF1ab > N > S > M = ORF10) were  
202 accumulating a higher number of synonymous mutations (dN/dS < 1). Interestingly, evaluation of

203 dN/dS at the level of individual genes showed higher ratios for the ORF3a and ORF8 genes  
204 (Figure 4A).

205 Significant purifying (negative) selection was observed in 12 out of the 35 codons evolving at  
206 dN/dS <1 using the FEL (12 sites) and FUBAR (1 site) methods; the codons were located in the  
207 ORF1ab, S, and N genes (Figure 4B). At these codons, increased fixation of synonymous  
208 substitutions seems to be favoring the phenotypic preservation of SARS-CoV-2 at specific  
209 residues of the proteins encoded by these genes. By tracking these mutations within different  
210 isolates, we observed that these changes could explain the divergence of different viruses within  
211 different phylogenetic groups. In some cases, mutations were associated with multiple isolates,  
212 supporting the relevance of these findings.

213 On the other hand, evidence of diversifying positive selection on non-synonymous sites was  
214 detected in just 4 of the 75 codons evolving at dN/dS >1 in genes ORF1ab and ORF3a, with the  
215 FUBAR and FEL methods providing the highest power of detection (Figure 4C). Based on this  
216 analysis, these four sites appear to be evolving under pervasive diversifying selection.

217 Interestingly, the detection of diversifying selection at codon 3606 (nsp6) was significantly  
218 supported by three different tests. Also, the selection of this site was observed in isolates from  
219 all three phylogenetic groups, thus supporting the reliability of this finding. However, the  
220 conservative nature of the amino acid substitution at this site (L-F; BSM62= 0) suggests this may  
221 not affect the phenotype of SARS-CoV-2. The same pattern was observed in codon 75 (D-E;  
222 BSM62= 2).

223 Conversely, based on the nature of the substitution (G-V; BSM62= -3), diversifying selection at  
224 codon 251 of the ORF3a gene may produce a biologically relevant effect on the phenotype of  
225 SARS-CoV-2. The same situation was observed at codon 2244 (I-T; BSM62= -1). Interestingly,  
226 change at codon position 251 is highly conserved within isolates of group C, suggesting that this  
227 change might have promoted the divergence of this group.

## 228 **Detection of directional selection**

229 To maximize the inference of potential sites experiencing positive selection, amino acid  
230 alignments of SARS-CoV-2 were analyzed using the DEPS algorithm. Overall, DEPS identified  
231 a total of 4 amino acid residues that are experiencing directional selection. Of these four  
232 residues, isoleucine (I) has the strongest bias, affecting 16 out of 19 sites evolving via directional  
233 selection (figure 5A).

234 The majority of selected sites were located in nonstructural proteins (nsp) encoded by the  
235 ORF1ab gene, with nsp3 accounting for the highest proportion (Figure 5B). Interestingly, a low  
236 proportion of the total number of predicted sites had a conservative amino acid substitution  
237 (residues at positions 902, 1769, 2235, and 2908), suggesting that the majority of substitutions  
238 may have an adaptive effect. In this context, it is remarkable that codon 84 of protein ORF8a is  
239 synapomorphic in all Group B sequences. Also, similar to previous algorithms, DEPS identified  
240 positive selection of residue 251 of ORF3a, supporting the potential significance of this site in  
241 the early evolution of SARS-CoV-2.

## 242 **Evidence of coevolution among sites**

243 Finally, we attempted to find coevolutionary correlations between different codons within the  
244 genome that result in the positive selection of sites. Analysis by BMG produced evidence of 14  
245 coevolving codon pairs; these interactions took place mostly within codons located within the  
246 ORF1ab gene (Figure 6). Although most of the interactions were detected between  
247 nonsynonymous codons, coevolution between codons 4090-4269 and 818-4320 was detected by  
248 a synonymous substitution at one of the codons. Also, based on the nature of the amino acid  
249 replacement, just 6 of the 14 interactions produced a potential biological significance in the  
250 replacement. Interestingly, 8 of the 14 interactions appeared associated with sites evolving under  
251 some type of positive selection, suggesting that the selection of these sites might be the result of  
252 epistatic events (Figure 7).

## 253 Discussion

254 Herein, we evaluated the phylogenetic and evolutionary dynamics of SARS-CoV-2 during the  
255 first months of the pandemic event in 2020. Our phylogenetic analysis revealed the complex  
256 dynamic of the spread of infection throughout the world, suggesting that multiple viral variants  
257 might have emerged in China before the start of the pandemic event. The evolutionary rate  
258 calculated for SARS-CoV-2 in this study was consistent not only with previous reports for  
259 SARS-CoV (Salemi et al., 2004, Zhao et al., 2004), but also with the rate for other RNA viruses  
260 (Sanjuan et al., 2010), explaining the high levels of identity at nucleotide and amino acid levels  
261 calculated for SARS-CoV-2 in our study. In this context, the high conservation observed in the  
262 genome of SARS-CoV-2 early during the pandemic might also be attributed to the unique RNA  
263 correction machinery of coronaviruses (Ferron et al., 2018).

264 However, and despite the relative genome stability observed in SARS-CoV-2 at this stage of the  
265 pandemic, we were able to describe the existence of 3 phylogenetic groups. Interestingly, our  
266 evolutionary analysis supported the hypothesis regarding early divergence events produced  
267 during the pandemic. By using a combination of different evolutionary algorithms, we detected  
268 multiple codon sites that may be promoting the divergence of SARS-CoV-2. However, two  
269 primary considerations must be addressed regarding the biological relevance of multiple sites  
270 detected in this study. First, a considerable number of polymorphisms were detected in just one  
271 viral isolate, which might be a consequence of the small number of viral isolates available at the  
272 time we started this study. Another possibility is that some of the polymorphisms might be due to  
273 sequencing errors. Additionally, some positively selected sites resulted in conservative  
274 replacements in terms of the similar nature of these amino acids, suggesting that they may not  
275 have a potential biological effect on the phenotype of SARS-CoV-2. However, the last  
276 assumption must be taken with caution, since experimental work on Chikungunya virus has  
277 shown that a conservative replacement of aspartic acid (D) with glutamic acid (E) at residue 350  
278 of the E1 protein is able to alter the affinity of a monoclonal antibody for this site (Tuekprakhon  
279 et al., 2018).

280 We are reporting three potentially relevant residues that may be driving the evolution of SARS-  
281 CoV-2 in human populations. Residue 251 of the ORF3a protein appears to be experiencing  
282 diversifying selection and might be related to the emergence of viruses in phylogenetic Group C.  
283 The early selection of this site might have a biological relevance since the ORF3a protein has  
284 been associated with virulence of human coronaviruses by controlling not only the expression of  
285 cytokines and chemokines but also inducing necrotic cell death (Shi et al., 2019).

286 Also, a residue located at position 84 of the ORF8a protein was found to be evolving under  
287 directional selection and might be related to the emergence of Group B. This protein has been  
288 implicated in viral pathogenesis by regulating the initial innate response (McBride and Fielding,  
289 2012, Shi et al., 2019). Interestingly, ORF3a and ORF8a accessory proteins had the highest  
290 dN/dS, indicating that a positive selective pressure is being exerted on these proteins during the  
291 pandemic.

292 Based on the known roles in tropism and virulence associated with the helicase of MERS-CoV,  
293 we consider it important to mention codon 5865 (ORF1ab) that was found to be evolving under  
294 directional selection and might be related to the divergence of 5 isolates from Washington, USA,  
295 forming a sub-cluster in Group B.

296 Finally, our analysis of coevolution revealed some potential epistatic interactions that might be  
297 driving the evolution of SARS-CoV-2. This mechanism has been proposed to explain the  
298 emergence of an Ebola virus variant in 2014 (Ibeh et al., 2016), and its relevance in the evolution  
299 of coronaviruses should be explored in future studies. Also, it is interesting to mention that most  
300 of the co-evolving sites were located in nsp3; given the role of this protein in the virulence of  
301 coronaviruses (Fehr et al., 2015), this observation may be key in understanding the evolution of  
302 SARS-CoV-2. Furthermore, since two of the interactions detected by BGM were associated with  
303 synonymous mutations, the relevance of this type of substitution to viral fitness should not be  
304 underestimated, since selection of synonymous substitutions has been reported in other RNA  
305 viruses like VSV (Novella et al., 2004, Velazquez-Salinas et al., 2018).

306 Collectively, our results describe the early evolutionary events of SARS-CoV-2 during the  
307 current pandemic and the findings may support the hypothesis that different variants of SARS-  
308 CoV-2 with disparate levels of virulence might be circulating in the world. This possibility  
309 might have an important impact on public health and measures to control the pandemic.  
310 Subsequent studies using reverse genetics will be needed to confirm the relevance of our findings  
311 connecting specific residue substitutions with different virus phenotypes.

312

### 313 **Author Contributions**

314 LV-S, and SZ conceived and designed the experiments. LV-S, SZ and SE performed the  
315 experiments. LV-S, SZ, SE, DG, IN, and MB analyzed the data. LV-S, SZ, SE, DG, IN, and  
316 MB wrote the manuscript.

### 317 **Funding**

318 This work was performed under USDA Research Service CRIS Project No. 8064-32000-060-  
319 00D.

320

### 321 **Conflict of Interest Statement**

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

### 324 **Acknowledgements**



325 We would like to particularly like to thank Melanie Prarat for editing the manuscript. Also, we  
326 thank Dr. Peter Simmonds for the use of the SSE program used for pairwise distance analysis.

327

## 328 **Figure Legends**

329 **Figure 1.** Sample summary . Description of the 86 SARS-Cov-2 full-length genome sequences  
330 included in this study. All sequences were obtained form I from the GenBank database,  
331 accession number, genome length, isolate name, source, host and country of origin are provided.  
332 N/A indicates information not available.

333 **Figure 2.** Phylogeny of SARS-Cov-2. A) A MCC tree reconstructed using 86 SARS-Cov-2 full-  
334 length genomes collected from patients naturally infected at different countries, showing the  
335 existence of 3 phylogenetic groups (A, B, C). Branches in the tree are colored based on the  
336 geographic location of each isolate. Presumed index case is pointed out with a red arrow. B)  
337 Average p-distance, and average nucleotide and amino acid differences between group A and  
338 groups B and C sequences, respectively. Analysis was conducted using the software MEGA 7.

339 **Figure 3.** Pairwise distance analysis. Pairwise distance analysis at: A) synonymous and B) non-  
340 synonymous nucleotide sites was conducted using the program Sequence Distances (software  
341 SSE). Red bars represent pairwise distance comparisons using an sliding window of 50  
342 nucleotides. Average nucleotide pairwise distance for different genes is shown at C)  
343 synonymous and D) non-synonymous sites. E) Fast-evolving synonymous and non-synonymous  
344 sites at each coding region are shown. For these sites evolutionary rates oscillated between 4.97  
345 and 4.95. Red numbers represent nucleotides at : 1)Leader protein, 2) nsp2, 3) nsp3, 4) nsp4, 5)  
346 nsp6, 6) nsp7, 7)nsp8, 8)nsp10, 9) RNA independent polymerase, 10) helicase, 11) 3'to 5'  
347 exonuclease, 12) endoRNase, 13) 2'-O-ribose methyltransferase.

348 **Figure 4.** Diversifying and purifying selection on SARS-CoV-2. A) General overview obtained  
349 by SLAC analysis, showing the evolutionary rate (dN-dS or dN/dS) along the genome and at  
350 individual genes of SARS-CoV-2. Statistically significant codons were inferred by multiple  
351 evolutionary tests used in this study. Red asterisks represent codons with significant evidence  
352 for selection, Codons evolving at B) Purifying (negative) or C) Diversifying (positive) selection  
353 are shown numbers in red represent evolutionary tests with significant values according to the  
354 analysis :SLAC, FEL, MEME (p-value=0.1) and FUBAR (posterior probability =0.9). The  
355 criteria for considering a site positively or negatively selected was based on their identification  
356 by at least one of the tests. The Phylogenetic group column (assigned according with figure 1A)  
357 shows also the isolates carrying the substitutions . Abbreviations: LP (Leader protein), 3LP (3C-  
358 like proteinase), n9 (nsp9), 3'-5' exo (3'-to 5' exonuclease), EN (endoRNase), and 2'M (2'-o-  
359 ribose methyltransferase).

360 **Figure 5.** Directional selection analysis on SARS-CoV-2. A)An amino acid alignment was  
361 evaluated by DEPS and 4 different residues producing 19 directionally evolving sites in the  
362 proteome of SARS-CoV-2 are reported. P-values show the statistical significance of each  
363 residue considering a model test of selection versus not selection. Bias term: Alignment-wide  
364 relative rate of substitution towards target residue. , Proportion of affected sites: Percentage of  
365 sites evolving under a directional model, versus a standard model with no directionality.  
366 Directionally evolving sites: Number of sites that show evidence of directional selection for focal

367 residue. B) Description of 19 directionally evolving sites. Sites were detected by Empirical  
368 Bayesian Factor (EBF) considering a cut-off of 100 or more. Numbers in red represent  
369 replacements between amino acids with different properties. The Phylogenetic group column  
370 (assigned according with figure 1A) shows also the isolates carrying the substitutions.

371 **Figure 6.** Coevolution between codon pairs in the genome of SARS-CoV-2. BMG analysis was  
372 conducted to detect coevolving codon pairs. Evidence of 14 coevolving codon pairs were  
373 detected and the specific locations of those in the genome of SARS-CoV-2 are presented.  
374 Posterior probability of pair associations was supported by Markov Chain Monte Carlo Analysis  
375 at cut-off of 50 or more. Numbers in red represent replacements between amino acids with  
376 different properties. The Phylogenetic group column (assigned according with figure 1A) shows  
377 also the isolates carrying the substitutions. \*<sup>1</sup> represents viral isolated where the changes were  
378 not detected.

379 **Figure 7.** Potential epistatic network in the genome of SARS-CoV-2. Coevolving pair codons  
380 involving sites evolving under positive selection at different regions in the genome of SARS-  
381 CoV-2 are shown with stars, suggesting that the selection of these sites may be result of epistasis.

382

## 383 References

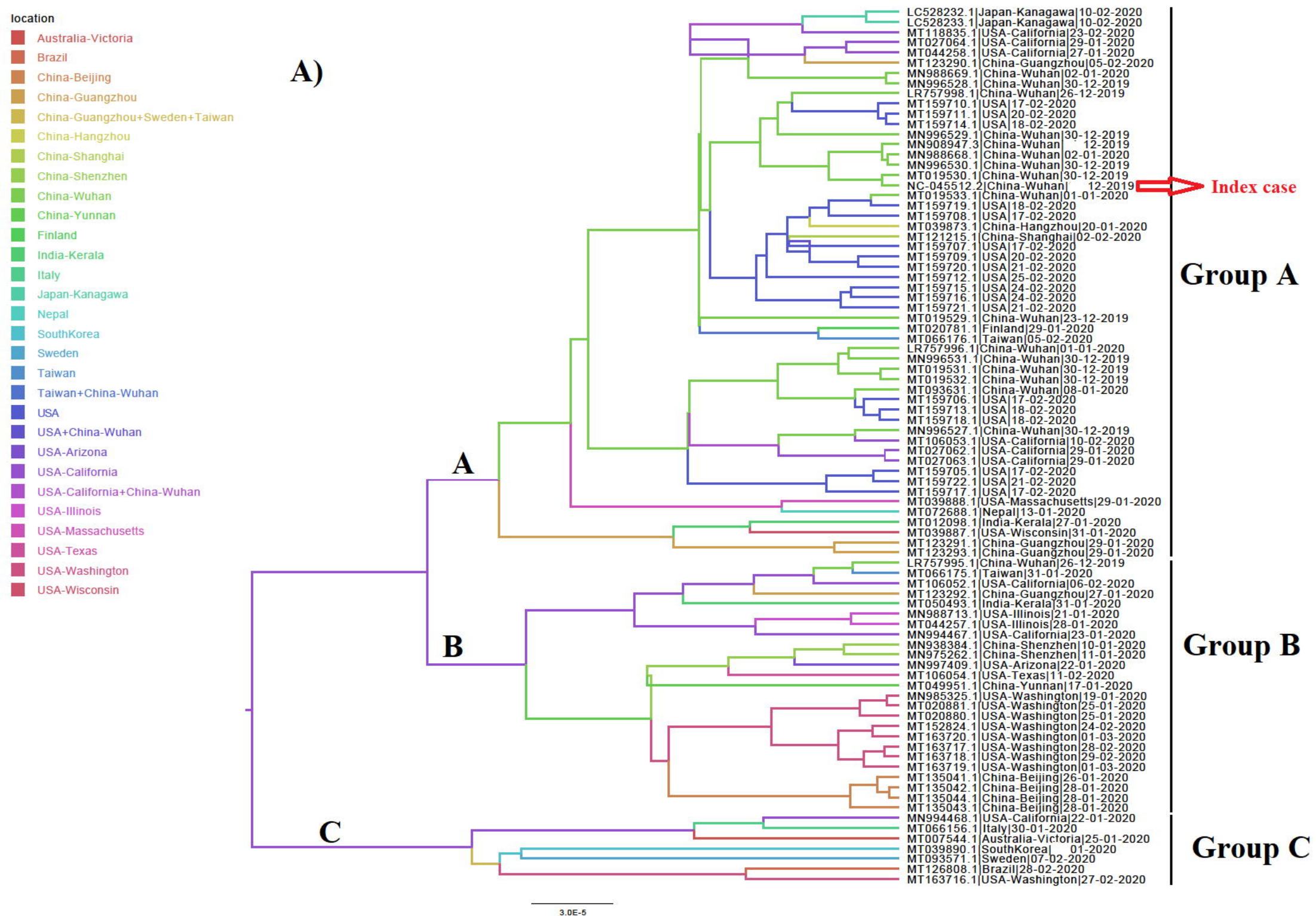
- 384  
385  
386 ANDERSEN, K. G., RAMBAUT, A., LIPKIN, W. I., HOLMES, E. C. & GARRY, R. F. 2020.  
387 The proximal origin of SARS-CoV-2. *Nature Medicine*.  
388 CORONAVIRIDAE STUDY GROUP OF THE INTERNATIONAL COMMITTEE ON  
389 TAXONOMY OF, V. 2020. The species Severe acute respiratory syndrome-related  
390 coronavirus: classifying 2019-nCoV and naming it SARS-CoV-2. *Nat Microbiol*.  
391 DELPORT, W., POON, A. F., FROST, S. D. & KOSAKOVSKY POND, S. L. 2010.  
392 Datamonkey 2010: a suite of phylogenetic analysis tools for evolutionary biology.  
393 *Bioinformatics*, 26, 2455-7.  
394 DRUMMOND, A. J., SUCHARD, M. A., XIE, D. & RAMBAUT, A. 2012. Bayesian  
395 phylogenetics with BEAUti and the BEAST 1.7. *Mol Biol Evol*, 29, 1969-73.  
396 FEHR, A. R., ATHMER, J., CHANNAPPANAVAR, R., PHILLIPS, J. M., MEYERHOLZ, D.  
397 K. & PERLMAN, S. 2015. The nsp3 macrodomain promotes virulence in mice with  
398 coronavirus-induced encephalitis. *J Virol*, 89, 1523-36.  
399 FERRON, F., SUBISSI, L., SILVEIRA DE MORAIS, A. T., LE, N. T. T., SEVAJOL, M.,  
400 GLUAIS, L., DECROLY, E., VONRHEIN, C., BRICOGNE, G., CANARD, B. &  
401 IMBERT, I. 2018. Structural and molecular basis of mismatch correction and ribavirin  
402 excision from coronavirus RNA. *Proc Natl Acad Sci U S A*, 115, E162-E171.  
403 FROST, S. D. W., MAGALIS, B. R. & KOSAKOVSKY POND, S. L. 2018. Neutral Theory and  
404 Rapidly Evolving Viral Pathogens. *Mol Biol Evol*, 35, 1348-1354.  
405 HENIKOFF, S. & HENIKOFF, J. G. 1992. Amino acid substitution matrices from protein  
406 blocks. *Proc Natl Acad Sci U S A*, 89, 10915-9.  
407 IBEH, N., NSHOGOZABAHIZI, J. C. & ARIS-BROSOUE, S. 2016. Both Epistasis and  
408 Diversifying Selection Drive the Structural Evolution of the Ebola Virus Glycoprotein  
409 Mucin-Like Domain. *J Virol*, 90, 5475-5484.

- 410 KOSAKOVSKY POND, S. L. & FROST, S. D. 2005. Not so different after all: a comparison of  
411 methods for detecting amino acid sites under selection. *Mol Biol Evol*, 22, 1208-22.
- 412 KOSAKOVSKY POND, S. L., POON, A. F., LEIGH BROWN, A. J. & FROST, S. D. 2008. A  
413 maximum likelihood method for detecting directional evolution in protein sequences and  
414 its application to influenza A virus. *Mol Biol Evol*, 25, 1809-24.
- 415 KOSAKOVSKY POND, S. L., POSADA, D., GRAVENOR, M. B., WOELK, C. H. & FROST,  
416 S. D. 2006. GARD: a genetic algorithm for recombination detection. *Bioinformatics*, 22,  
417 3096-8.
- 418 KUMAR, S., STECHER, G. & TAMURA, K. 2016. MEGA7: Molecular Evolutionary Genetics  
419 Analysis Version 7.0 for Bigger Datasets. *Mol Biol Evol*, 33, 1870-4.
- 420 MCBRIDE, R. & FIELDING, B. C. 2012. The role of severe acute respiratory syndrome  
421 (SARS)-coronavirus accessory proteins in virus pathogenesis. *Viruses*, 4, 2902-23.
- 422 MURRELL, B., MOOLA, S., MABONA, A., WEIGHILL, T., SHEWARD, D.,  
423 KOSAKOVSKY POND, S. L. & SCHEFFLER, K. 2013. FUBAR: a fast, unconstrained  
424 bayesian approximation for inferring selection. *Mol Biol Evol*, 30, 1196-205.
- 425 MURRELL, B., WERTHEIM, J. O., MOOLA, S., WEIGHILL, T., SCHEFFLER, K. &  
426 KOSAKOVSKY POND, S. L. 2012. Detecting individual sites subject to episodic  
427 diversifying selection. *PLoS Genet*, 8, e1002764.
- 428 NOVELLA, I. S., ZARATE, S., METZGAR, D. & EBENDICK-CORPUS, B. E. 2004. Positive  
429 selection of synonymous mutations in vesicular stomatitis virus. *J Mol Biol*, 342, 1415-  
430 21.
- 431 PHAN, M. V. T., NGO TRI, T., HONG ANH, P., BAKER, S., KELLAM, P. & COTTEN, M.  
432 2018. Identification and characterization of Coronaviridae genomes from Vietnamese  
433 bats and rats based on conserved protein domains. *Virus Evol*, 4, vey035.
- 434 POON, A. F., LEWIS, F. I., POND, S. L. & FROST, S. D. 2007. An evolutionary-network  
435 model reveals stratified interactions in the V3 loop of the HIV-1 envelope. *PLoS Comput  
436 Biol*, 3, e231.
- 437 RAMBAUT, A., DRUMMOND, A. J., XIE, D., BAELE, G. & SUCHARD, M. A. 2018.  
438 Posterior Summarization in Bayesian Phylogenetics Using Tracer 1.7. *Syst Biol*, 67, 901-  
439 904.
- 440 SALEMI, M., FITCH, W. M., CICCOCCHI, M., RUIZ-ALVAREZ, M. J., REZZA, G. & LEWIS,  
441 M. J. 2004. Severe acute respiratory syndrome coronavirus sequence characteristics and  
442 evolutionary rate estimate from maximum likelihood analysis. *J Virol*, 78, 1602-3.
- 443 SANJUAN, R., NEBOT, M. R., CHIRICO, N., MANSKY, L. M. & BELSHAW, R. 2010. Viral  
444 mutation rates. *J Virol*, 84, 9733-48.
- 445 SHI, C. S., NABAR, N. R., HUANG, N. N. & KEHRL, J. H. 2019. SARS-Coronavirus Open  
446 Reading Frame-8b triggers intracellular stress pathways and activates NLRP3  
447 inflammasomes. *Cell Death Discov*, 5, 101.
- 448 SIMMONDS, P. 2012. SSE: a nucleotide and amino acid sequence analysis platform. *BMC Res  
449 Notes*, 5, 50.
- 450 SUN, P., QIE, S., LIU, Z., REN, J., LI, K. & XI, J. 2020. Clinical characteristics of hospitalized  
451 patients with SARS-CoV-2 infection: A single arm meta-analysis. *J Med Virol*.
- 452 TUEKPRAXHON, A., NAKAYAMA, E. E., BARTHOLOMEEUSEN, K., PUIPROM, O.,  
453 SASAKI, T., HUIITS, R., LUPLERTLOP, N., KOSOLTANAPIWAT, N., MANEEKAN,  
454 P., ARIEN, K. K., SHIODA, T. & LEAUNGWUTIWONG, P. 2018. Variation at

- 455 position 350 in the Chikungunya virus 6K-E1 protein determines the sensitivity of  
456 detection in a rapid E1-antigen test. *Sci Rep*, 8, 1094.
- 457 VELAZQUEZ-SALINAS, L., PAUSZEK, S. J., STENFELDT, C., O'HEARN, E. S.,  
458 PACHECO, J. M., BORCA, M. V., VERDUGO-RODRIGUEZ, A., ARZT, J. &  
459 RODRIGUEZ, L. L. 2018. Increased Virulence of an Epidemic Strain of Vesicular  
460 Stomatitis Virus Is Associated With Interference of the Innate Response in Pigs. *Front*  
461 *Microbiol*, 9, 1891.
- 462 WAN, Y., SHANG, J., GRAHAM, R., BARIC, R. S. & LI, F. 2020. Receptor Recognition by  
463 the Novel Coronavirus from Wuhan: an Analysis Based on Decade-Long Structural  
464 Studies of SARS Coronavirus. *J Virol*, 94.
- 465 WEAVER, S., SHANK, S. D., SPIELMAN, S. J., LI, M., MUSE, S. V. & KOSAKOVSKY  
466 POND, S. L. 2018. Datamonkey 2.0: A Modern Web Application for Characterizing  
467 Selective and Other Evolutionary Processes. *Mol Biol Evol*, 35, 773-777.
- 468 WEISS, S. R. & NAVAS-MARTIN, S. 2005. Coronavirus pathogenesis and the emerging  
469 pathogen severe acute respiratory syndrome coronavirus. *Microbiol Mol Biol Rev*, 69,  
470 635-64.
- 471 WU, F., ZHAO, S., YU, B., CHEN, Y. M., WANG, W., SONG, Z. G., HU, Y., TAO, Z. W.,  
472 TIAN, J. H., PEI, Y. Y., YUAN, M. L., ZHANG, Y. L., DAI, F. H., LIU, Y., WANG, Q.  
473 M., ZHENG, J. J., XU, L., HOLMES, E. C. & ZHANG, Y. Z. 2020. A new coronavirus  
474 associated with human respiratory disease in China. *Nature*, 579, 265-269.
- 475 ZHAO, Z., LI, H., WU, X., ZHONG, Y., ZHANG, K., ZHANG, Y. P., BOERWINKLE, E. &  
476 FU, Y. X. 2004. Moderate mutation rate in the SARS coronavirus genome and its  
477 implications. *BMC Evol Biol*, 4, 21.
- 478

Genbank accession number	Genome bases	isolation source	Host	Country	State/City/Province	Collection Date
LC528232.1	29902	throat swab	Homo sapiens	Japan	Kanagawa	10-Feb-20
LC528233.1	29902	throat swab	Homo sapiens	Japan	Kanagawa	10-Feb-20
LR757995.1	29872	N/A	Homo sapiens	China	Wuhan	26-Dec-19
LR757996.1	29868	N/A	Homo sapiens	China	Wuhan	1-Jan-20
LR757998.1	29866	N/A	Homo sapiens	China	Wuhan	26-Dec-19
MN908947.3	29903	N/A	Homo sapiens	China	Wuhan	Dec-19
MN938384.1	29838	nasopharyngeal swab	Homo sapiens	China	Shenzhen	10-Jan-20
MN975262.1	29891	sputum	Homo sapiens	China	Shenzhen	11-Jan-20
MN985325.1	29882	oropharyngeal swab	Homo sapiens	USA	Washington	19-Jan-20
MN988668.1	29881	bronchoalveolar lavage fluid	Homo sapiens	China	Wuhan	2-Jan-20
MN988669.1	29881	bronchoalveolar lavage fluid	Homo sapiens	China	Wuhan	2-Jan-20
MN988713.1	29882	sputum	Homo sapiens	USA	Illinois	21-Jan-20
MN994467.1	29882	nasopharyngeal swab	Homo sapiens	USA	California	23-Jan-20
MN994468.1	29883	nasopharyngeal swab	Homo sapiens	USA	California	22-Jan-20
MN996527.1	29825	bronchoalveolar lavage fluid	Homo sapiens	China	Wuhan	30-Dec-19
MN996528.1	29891	bronchoalveolar lavage fluid	Homo sapiens	China	Wuhan	30-Dec-19
MN996529.1	29852	bronchoalveolar lavage fluid	Homo sapiens	China	Wuhan	30-Dec-19
MN996530.1	29854	bronchoalveolar lavage fluid	Homo sapiens	China	Wuhan	30-Dec-19
MN996531.1	29857	bronchoalveolar lavage fluid	Homo sapiens	China	Wuhan	30-Dec-19
MN997409.1	29882	buccal swab	Homo sapiens	USA	Arizona	22-Jan-20
MT007544.1	29893	N/A	Homo sapiens	Australia	Victoria	25-Jan-20
MT012098.1	29854	throat swab	Homo sapiens	India	Kerala	27-Jan-20
MT019529.1	29899	bronchoalveolar lavage fluid	Homo sapiens; male; age 65	China	Wuhan	23-Dec-19
MT019530.1	29889	bronchoalveolar lavage fluid	Homo sapiens; female; age 49	China	Wuhan	30-Dec-19
MT019531.1	29899	bronchoalveolar lavage fluid	Homo sapiens; male; age 41	China	Wuhan	30-Dec-19
MT019532.1	29890	bronchoalveolar lavage fluid	Homo sapiens; female; age 52	China	Wuhan	30-Dec-19
MT019533.1	29883	bronchoalveolar lavage fluid	Homo sapiens; male; age 61	China	Wuhan	1-Jan-20
MT020781.2	29806	N/A	Homo sapiens	Finland	N/A	29-Jan-20
MT020880.1	29882	nasopharyngeal swab	Homo sapiens	USA	Washington	25-Jan-20
MT020881.1	29882	oropharyngeal swab	Homo sapiens	USA	Washington	25-Jan-20
MT027062.1	29882	nasopharyngeal swab	Homo sapiens	USA	California	29-Jan-20
MT027063.1	29882	oropharyngeal swab	Homo sapiens	USA	California	29-Jan-20
MT027064.1	29882	oropharyngeal swab	Homo sapiens	USA	California	29-Jan-20
MT039873.1	29833	sputum	Homo sapiens; male	China	Hangzhou	20-Jan-20
MT039887.1	29879	nasopharyngeal swab	Homo sapiens	USA	Wisconsin	31-Jan-20
MT039888.1	29882	oropharyngeal swab	Homo sapiens	USA	Massachusetts	29-Jan-20
MT039890.1	29903	N/A	Homo sapiens	South Korea	N/A	Jan-20
MT044257.1	29882	sputum	Homo sapiens	USA	Illinois	28-Jan-20
MT044258.1	29858	respiratory swab	Homo sapiens	USA	California	27-Jan-20
MT049951.1	29903	sputum	Homo sapiens	China	Yunnan	17-Jan-20
MT050493.1	29851	throat swab	Homo sapiens	India	Kerala	31-Jan-20
MT066156.1	29867	sputum	Homo sapiens	Italy	N/A	30-Jan-20
MT066175.1	29870	N/A	Homo sapiens	Taiwan	N/A	31-Jan-20
MT066176.1	29870	N/A	Homo sapiens	Taiwan	N/A	5-Feb-20
MT072688.1	29811	oropharyngeal swab	Homo sapiens	Nepal	N/A	13-Jan-20
MT093571.1	29886	N/A	Homo sapiens	Sweden	N/A	7-Feb-20
MT093631.2	29860	throat swab	Homo sapiens	China	Wuhan	8-Jan-20
MT106052.1	29882	nasopharyngeal swab	Homo sapiens	USA	California	6-Feb-20
MT106053.1	29882	nasopharyngeal swab	Homo sapiens	USA	California	10-Feb-20
MT106054.1	29882	sputum	Homo sapiens	USA	Texas	11-Feb-20
MT118835.1	29882	bronchoalveolar lavage	Homo sapiens	USA	California	23-Feb-20
MT121215.1	29945	throat swab	Homo sapiens	China	Shanghai	2-Feb-20
MT123290.1	29891	oropharyngeal swab	Homo sapiens	China	Guangzhou	5-Feb-20
MT123291.2	29982	bronchoalveolar lavage fluid	Homo sapiens	China	Guangzhou	29-Jan-20
MT123292.2	29923	sputum	Homo sapiens	China	Guangzhou	27-Jan-20
MT123293.2	29871	stool	Homo sapiens	China	Guangzhou	29-Jan-20
MT126808.1	29876	nasopharyngeal swab	Homo sapiens	Brazil	N/A	28-Feb-20
MT135041.1	29903	N/A	Homo sapiens	China	Beijing	26-Jan-20
MT135042.1	29903	N/A	Homo sapiens	China	Beijing	28-Jan-20
MT135043.1	29903	N/A	Homo sapiens	China	Beijing	28-Jan-20
MT135044.1	29903	N/A	Homo sapiens	China	Beijing	28-Jan-20
MT152824.1	29878	mid-nasal swab	Homo sapiens	USA	Washington	24-Feb-20
MT159705.1	29882	nasopharyngeal swab	Homo sapiens	USA	N/A	17-Feb-20
MT159706.1	29882	nasopharyngeal swab	Homo sapiens	USA	N/A	17-Feb-20
MT159707.1	29882	nasopharyngeal swab	Homo sapiens	USA	N/A	17-Feb-20
MT159708.1	29882	nasopharyngeal swab	Homo sapiens	USA	N/A	17-Feb-20
MT159709.1	29882	oropharyngeal swab	Homo sapiens	USA	N/A	20-Feb-20
MT159710.1	29882	nasopharyngeal swab	Homo sapiens	USA	N/A	17-Feb-20
MT159711.1	29882	nasopharyngeal swab	Homo sapiens	USA	N/A	20-Feb-20
MT159712.1	29882	oropharyngeal swab	Homo sapiens	USA	N/A	25-Feb-20
MT159713.1	29882	nasopharyngeal swab	Homo sapiens	USA	N/A	18-Feb-20
MT159714.1	29882	nasopharyngeal swab	Homo sapiens	USA	USA	18-Feb-20
MT159715.1	29882	nasopharyngeal swab	Homo sapiens	USA	N/A	24-Feb-20
MT159716.1	29687	nasopharyngeal swab	Homo sapiens	USA	N/A	24-Feb-20
MT159717.1	29882	nasopharyngeal swab	Homo sapiens	USA	N/A	17-Feb-20
MT159718.1	29882	nasopharyngeal swab	Homo sapiens	USA	N/A	18-Feb-20
MT159719.1	29882	nasopharyngeal swab	Homo sapiens	USA	N/A	18-Feb-20
MT159720.1	29882	nasopharyngeal swab	Homo sapiens	USA	N/A	21-Feb-20
MT159721.1	29882	oropharyngeal swab	Homo sapiens	USA	N/A	21-Feb-20
MT159722.1	29882	nasopharyngeal swab	Homo sapiens	USA	N/A	21-Feb-20
MT163716.1	29903	N/A	Homo sapiens	USA	Washington	27-Feb-20
MT163717.1	29897	N/A	Homo sapiens	USA	Washington	28-Feb-20
MT163718.1	29903	N/A	Homo sapiens	USA	Washington	29-Feb-20
MT163719.1	29903	N/A	Homo sapiens	USA	Washington	1-Mar-20
MT163720.1	29732	N/A	Homo sapiens	USA	Washington	1-Mar-20
NC_045512.2	29903	N/A	Homo sapiens	China	Wuhan	Dec-19

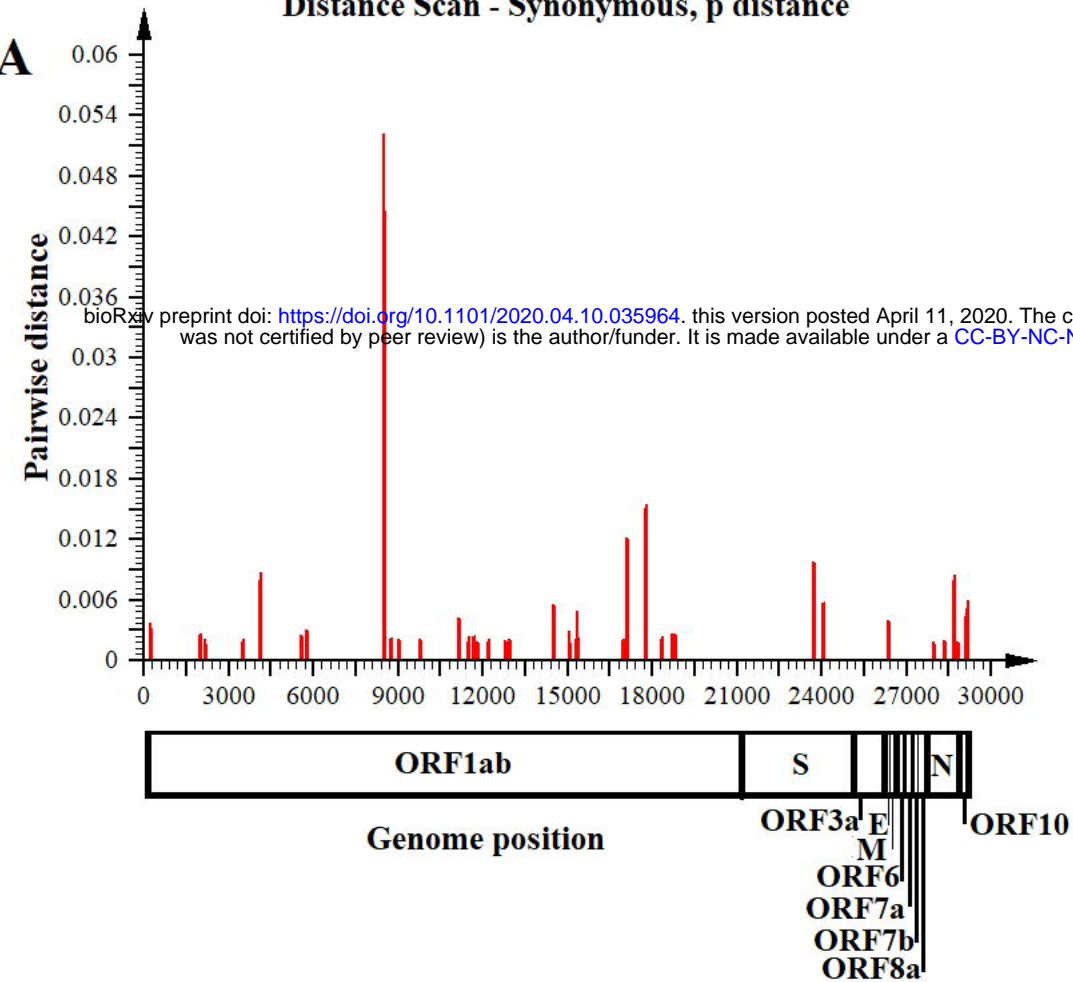
bioRxiv preprint doi: <https://doi.org/10.1101/2020.04.10.035964>; this version posted April 11, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. It is made available under aCC-BY-NC-ND 4.0 International license.



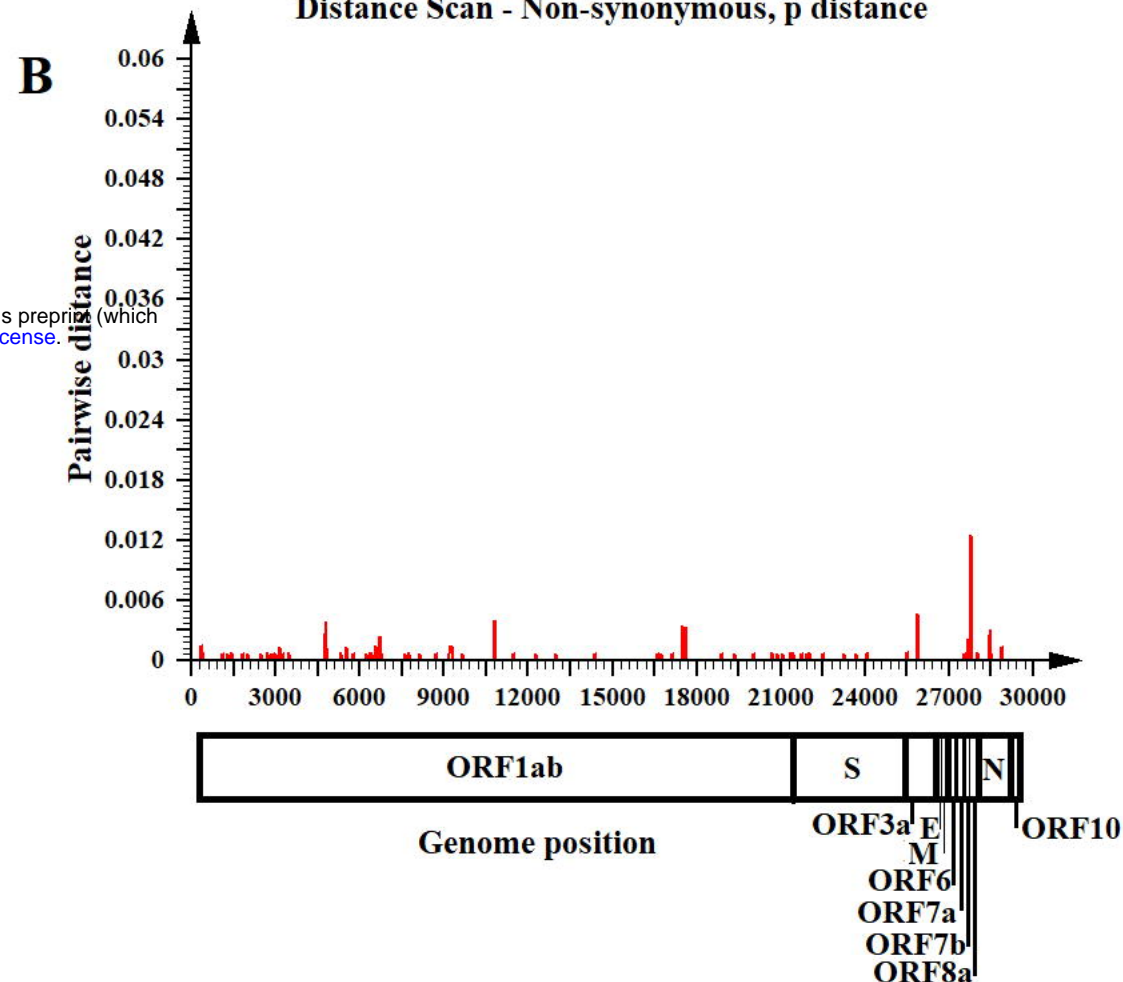
**B)**

	Average p-distance	Average nucleotide differences	Average amino acid differences
Group B	0.000128	3.75	2
Group C	0.000191	5.78	3.67

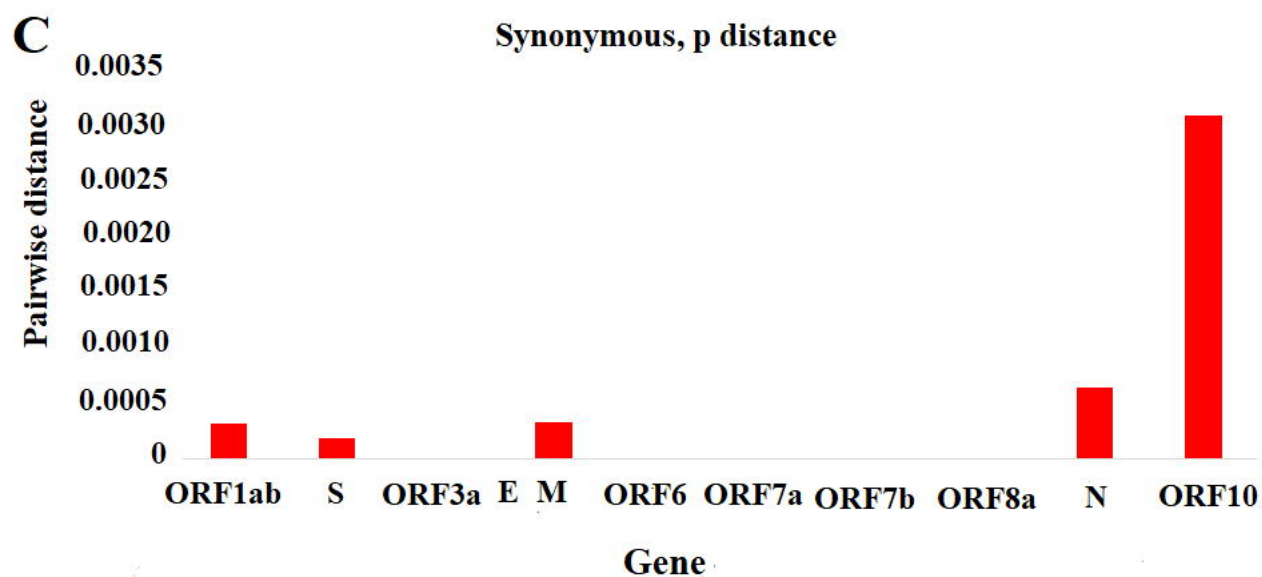
**A** Distance Scan - Synonymous, p distance



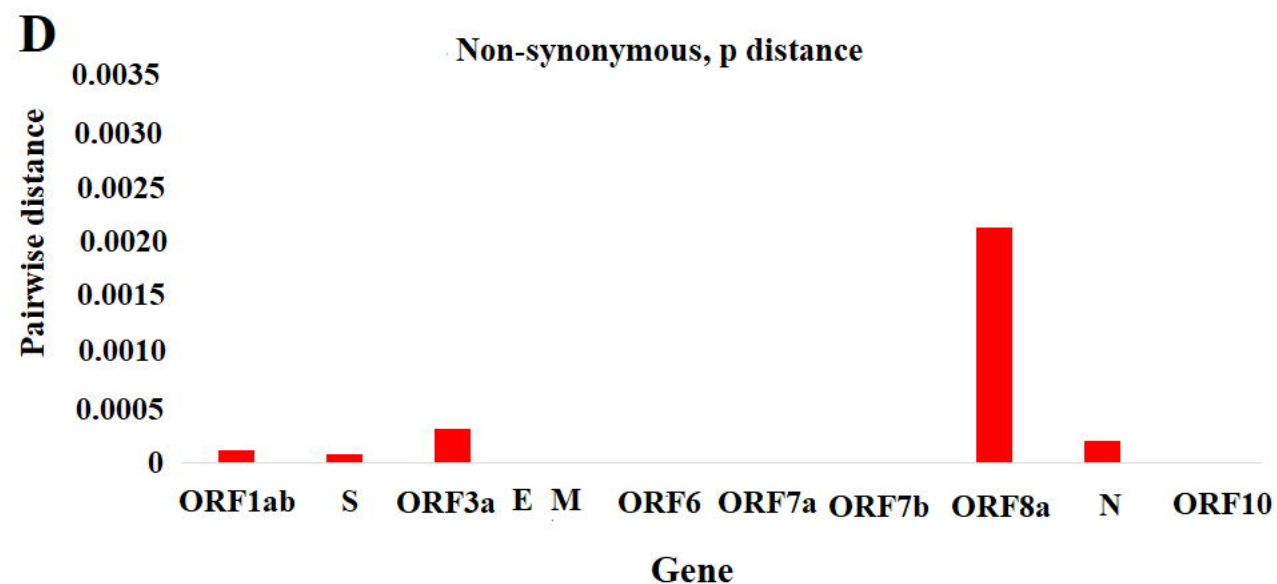
**B** Distance Scan - Non-synonymous, p distance



**C** Synonymous, p distance

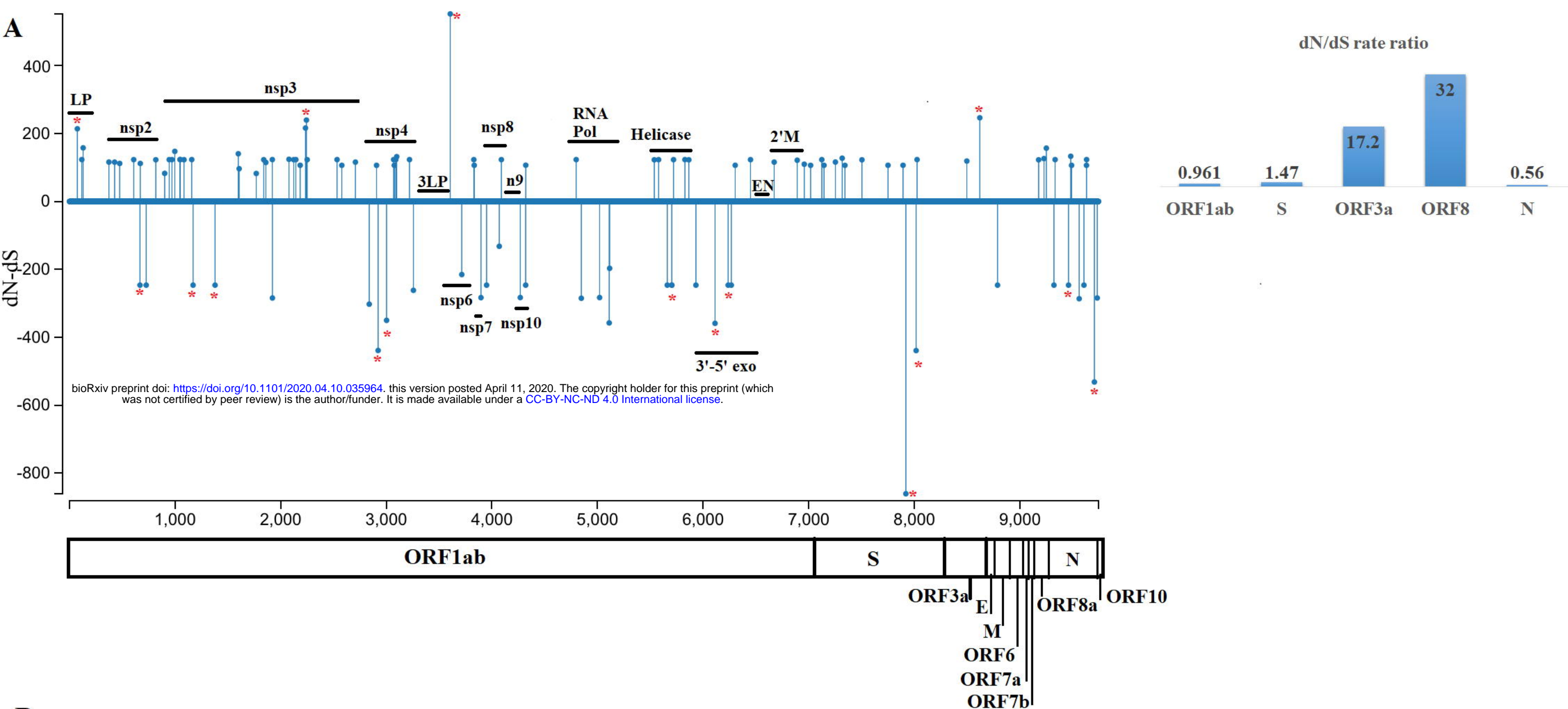


**D** Non-synonymous, p distance



**E**

Gene	Total of evolving sites >1	# of synonymous sites	Gene nucleotide position	# of non-synonymous sites	Gene nucleotide position
ORF1ab	82	28	(349) <sup>1</sup> , (2004, 2181) <sup>2</sup> , (3513, 4137, 5766) <sup>3</sup> , (8517, 8769, 9009, 9771) <sup>4</sup> , (11145) <sup>5</sup> , (11691) <sup>6</sup> , (11850, 12208) <sup>7</sup> , (12807, 12960) <sup>8</sup> , (14541, 15060, 15333, 15343) <sup>9</sup> , (16983, 17109, 17112) <sup>10</sup> , (17796, 18339, 18711, 18801) <sup>11</sup>	54	(225, 389) <sup>1</sup> , (1120, 1283, 1426, 1826, 2012, 2452) <sup>2</sup> , (2706, 2834, 2912, 2994, 3141, 3146, 3253, 3473, 4797, 4819, 5307, 5519, 5580, 5761, 6236, 6371, 6430, 6554, 6703, 6731, 6751, 7601, 7736, 8123) <sup>3</sup> , (8722, 9226, 9269, 9659) <sup>4</sup> , (10818, 11485, 11499) <sup>5</sup> , (12269) <sup>7</sup> , (12961) <sup>8</sup> , (14393) <sup>9</sup> , (16613, 16736, 17159, 17483, 17594) <sup>10</sup> , (18911, 19346) <sup>11</sup> , (20017) <sup>12</sup> , (20672, 20878, 21052) <sup>13</sup>
S	11	3	471, 2472, 2763	8	82, 145, 662, 741, 1223, 1963, 2390, 2789
ORF3a	2	0		2	383, 752
M	1	1	207	0	
ORF8a	3	0		3	32, 184, 251
N	10	5	105, 519, 822, 957	5	136, 581, 605, 1028, 1030
ORF10	1	1	78	0	



**C**

Gene	Gene codon position	SLAC		FEL		MEME		FUBAR		Inferred substitution	Phylogenetic group
		dN-dS	P [dN/dS < 1]	dN/dS	p-value	dN/neutral evolution component	p-value	dN-dS	Prob[dS<dN]		
ORF1ab/leader protein	75	214.12	0.60	Infinity	0.25	293.63	0.26	31.52	<b>0.94</b>	(D) GAU <sub>84(D)</sub> GAA <sub>2(E)</sub>	Group B (MN988713.1, MT04427.1)
ORF1ab/nsp3	2244	239.71	0.47	Infinity	<b>0.10</b>	480.91	0.12	31.52	<b>0.95</b>	(I) AUC <sub>84(I)</sub> ACC <sub>2(T)</sub>	Group A (MT019531.1, MT123293.1)
ORF1ab/nsp6	3606	552.28	0.14	Infinity	<b>0.08</b>	3129.69	<b>0.10</b>	41.36	<b>0.97</b>	(U) UUG <sub>80(U)</sub> UUU <sub>5(F)</sub> UUC <sub>1*</sub>	Group A (LC528232.1, LC5282233.1) Group B (MT049951.1*) Group C (MT163716.1, MT049951.1)
ORF3a	251	246.55	0.44	Infinity	0.19	517.70	0.21	31.38	<b>0.95</b>	(G) GGU <sub>79(G)</sub> GUU <sub>7(V)</sub>	Group C (all 7 taxon)



**A**

Residue	p-value	Bias term	Proportion of affected sites	Directionally evolving sites
C	0.0028	10.89	50.00%	1
I	0	716.43	0.69%	16
S	0	47.28	1.71%	1
V	0	24.36	3.80%	1

**B**

Protein	Site	Composition in the alignment	Reconstructed most common ancestor at site	Inferred Substitutions	DEPS EBF	Blosum 62 Matrix Score	Phylogenetic group
ORF1ab/nsp2	609	T <sub>85</sub> I <sub>1</sub>	T	I <sub>1</sub> ↔ <sub>0</sub> T	I: 375.3	<b>-1</b>	Group A (MT027064.1)
ORF1ab/nsp3	902	M <sub>85</sub> I <sub>1</sub>	M	I <sub>1</sub> ↔ <sub>0</sub> M	I: 224.4	1	Group C (MT039890.1)
ORF1ab/nsp3	945	T <sub>85</sub> I <sub>1</sub>	T	I <sub>1</sub> ↔ <sub>0</sub> T	I: 374.2	<b>-1</b>	Group A (MT159717.1)
ORF1ab/nsp3	1769	M <sub>85</sub> I <sub>1</sub>	M	I <sub>1</sub> ↔ <sub>0</sub> M	I: 223.3	1	Group C (MT163716.1)
ORF1ab/nsp3	1840	T <sub>84</sub> I <sub>2</sub>	T	I <sub>2</sub> ↔ <sub>0</sub> T	I:>10 <sup>5</sup>	<b>-1</b>	Group B (MT152824.1, MT163720.1)
ORF1ab/nsp3	2124	T <sub>85</sub> I <sub>1</sub>	T	I <sub>1</sub> ↔ <sub>0</sub> T	I: 375.3	<b>-1</b>	Group A (MT159712.1)
ORF1ab/nsp3	2185	S <sub>84</sub> I <sub>2</sub>	S	I <sub>2</sub> ↔ <sub>0</sub> S	I:>10 <sup>5</sup>	<b>-2</b>	Group A (MT123291.1, MT123293.1)
ORF1ab/nsp3	2235	L <sub>85</sub> I <sub>1</sub>	L	I <sub>1</sub> ↔ <sub>0</sub> L	I: 367.5	2	Group A (MT019529.1)
ORF1ab/nsp4	2908	F <sub>85</sub> I <sub>1</sub>	F	F <sub>0</sub> ↔ <sub>1</sub> I	I: 558.4	1	Group A (MN996531.1)
ORF1ab/nsp4	3090	T <sub>85</sub> I <sub>1</sub>	T	I <sub>1</sub> ↔ <sub>0</sub> T	I: 375.3	<b>-1</b>	Group A (LR757998.1)
ORF1ab/nsp8	4090	T <sub>85</sub> I <sub>1</sub>	T	I <sub>1</sub> ↔ <sub>0</sub> T	I: 375.3	<b>-1</b>	Group A (MT123292.1)
ORF1ab/helicase	5538	T <sub>85</sub> I <sub>1</sub>	T	I <sub>1</sub> ↔ <sub>0</sub> T	I: 378.5	<b>-1</b>	Group B (MT050493.1)
ORF1ab/helicase	5579	T <sub>85</sub> I <sub>1</sub>	T	I <sub>1</sub> ↔ <sub>0</sub> T	I: 375.3	<b>-1</b>	Group C (MN994468.1)
ORF1ab/helicase	5865	Y <sub>81</sub> C <sub>5</sub>	Y	C <sub>5</sub> ↔ <sub>0</sub> Y	C: 695.1	<b>-2</b>	Group B (MT163717.1, MT163718.1, MT163719.1, MT163720.1, MT152824.1)
ORF1ab/ 3'-5' exonuclease	6449	T <sub>85</sub> I <sub>1</sub>	T	I <sub>1</sub> ↔ <sub>0</sub> T	I: 375.5	<b>-1</b>	Group A (MT123291.1)
S	408	R <sub>85</sub> I <sub>1</sub>	R	I <sub>1</sub> ↔ <sub>0</sub> R	I: 652.7	<b>-3</b>	Group A MT012098.1
ORF3a	251	G <sub>79</sub> V <sub>7</sub>	G	G <sub>0</sub> ↔ <sub>7</sub> V	V:>10 <sup>5</sup>	<b>-3</b>	Group C (all 7 taxon)
ORF8a	11	T <sub>85</sub> I <sub>1</sub>	T	I <sub>1</sub> ↔ <sub>0</sub> T	I: 376.3	<b>-1</b>	Group B (MT106054.1)
ORF8a	84	L <sub>61</sub> S <sub>25</sub>	L	L <sub>0</sub> ↔ <sub>25</sub> S	S:>10 <sup>5</sup>	<b>-2</b>	Group B (all 25 taxon)

Gene site 1	Codon position	Gene site 2	Codon position	Probability that sites 1 and 2 are not conditionally independent	Inferred substitution site 1	Blosum 62 matrix score	Inferred substitution site 2	Blosum 62 matrix score	Phylogenetic group
ORF1ab/leader protein	75	ORF8a	62	0.5	<sub>(D)</sub> GAU <sub>84</sub> <sub>(E)</sub> GAA <sub>2</sub>	2	<sub>(V)</sub> GUG <sub>83</sub> <sub>(L)</sub> CUG <sub>3</sub>	1	Group B (MT044257.1, MN988713.1) MN994467.1* <sup>1</sup>
ORF1ab/leader protein	117	ORF1ab/nsp3	1607	0.76	<sub>(A)</sub> GCU <sub>84</sub> <sub>(T)</sub> ACU <sub>2</sub>	0	<sub>(I)</sub> AUA <sub>84</sub> <sub>(V)</sub> GUA <sub>2</sub>	3	Group A (MT027062.1, MT027063.1)
ORF1ab/leader protein	130	ORF1ab/nsp3	2244	0.56	<sub>(G)</sub> GGA <sub>85</sub> <sub>(E)</sub> GAA <sub>1</sub>	<b>-2</b>	<sub>(I)</sub> AUC <sub>83</sub> <sub>(T)</sub> ACC <sub>2</sub>	<b>-1</b>	Group A (MT123293.1) MT019353.1* <sup>1</sup>
ORF1ab/nsp2	609	S	49	0.74	<sub>(T)</sub> ACU <sub>85</sub> <sub>(I)</sub> AUU <sub>1</sub>	<b>-1</b>	<sub>(H)</sub> CUA <sub>85</sub> <sub>(Y)</sub> UAU <sub>1</sub>	2	Group A (MT027064.1)
ORF1ab/nsp2	818	ORF1ab/nsp10	4320	0.76	<sub>(G)</sub> GGU <sub>85</sub> <sub>(S)</sub> AGU <sub>1</sub>	0	<sub>(S)</sub> UCC <sub>85</sub> <sub>(S)</sub> UCG <sub>1</sub>	N/A	Group C (MT093571.1)
ORF1ab/nsp3	1047	S	655	0.73	<sub>(E)</sub> GAA <sub>85</sub> <sub>(D)</sub> GAC <sub>1</sub>	2	<sub>(H)</sub> CAU <sub>85</sub> <sub>(Y)</sub> UAU <sub>1</sub>	2	Group B (MT163720.1)
ORF1ab/nsp3	1049	ORF1ab/nsp3	1769	0.76	<sub>(A)</sub> GCU <sub>85</sub> <sub>(V)</sub> GUU <sub>1</sub>	0	<sub>(M)</sub> AUG <sub>85</sub> <sub>(I)</sub> AUU <sub>1</sub>	1	Group C (MT163716.1)
ORF1ab/nsp3	2124	ORF1ab/nsp6	3829	0.75	<sub>(T)</sub> ACU <sub>85</sub> <sub>(I)</sub> AUU <sub>1</sub>	<b>-1</b>	<sub>(L)</sub> CUC <sub>85</sub> <sub>(I)</sub> UUC <sub>1</sub>	2	Group A (MT159712.1)
ORF1ab/nsp3	2235	ORF1ab/nsp6	3833	0.74	<sub>(L)</sub> CUA <sub>85</sub> <sub>(I)</sub> AUA <sub>1</sub>	2	<sub>(N)</sub> AAU <sub>85</sub> <sub>(K)</sub> AAA <sub>1</sub>	0	Group B (LR757998.1)
ORF1ab/nsp3	2251	ORF1ab/2'-O-ribose methyltransferase	6958	0.74	<sub>(G)</sub> GGU <sub>85</sub> <sub>(S)</sub> AGU <sub>1</sub>	0	<sub>(K)</sub> AAG <sub>85</sub> <sub>(R)</sub> AGG <sub>1</sub>	2	Group A (MN99629.1)
ORF1ab/nsp3	2579	ORF1ab/nsp4	3090	0.72	<sub>(D)</sub> GAU <sub>85</sub> <sub>(A)</sub> GCU <sub>1</sub>	<b>-2</b>	<sub>(T)</sub> ACU <sub>85</sub> <sub>(I)</sub> AUU <sub>1</sub>	<b>-1</b>	Group A (MN996531.1)
ORF1ab/nsp3	2708	ORF1ab/nsp4	2908	0.73	<sub>(N)</sub> AAC <sub>85</sub> <sub>(S)</sub> AGC <sub>1</sub>	1	<sub>(F)</sub> UUU <sub>85</sub> <sub>(I)</sub> AUU <sub>1</sub>	0	Group A (MT019529.1)
ORF1ab/nsp8	4090	ORF1ab/nsp10	4269	0.73	<sub>(T)</sub> ACU <sub>85</sub> <sub>(I)</sub> AUU <sub>1</sub>	<b>-1</b>	<sub>(F)</sub> UUC <sub>85</sub> <sub>(F)</sub> UUU <sub>1</sub>	N/A	Group B (MT123292.1)
ORF1ab/3'-to-5' exonuclease	6304	ORF8a	11	0.73	<sub>(D)</sub> GAU <sub>85</sub> <sub>(A)</sub> GCU <sub>1</sub>	<b>-2</b>	<sub>(T)</sub> ACA <sub>85</sub> <sub>(I)</sub> AUA <sub>1</sub>	<b>-1</b>	Group B (MT106054.1)

