

1 **Title:** SARS-CoV-2 amino acid substitutions widely spread in the human population are
2 mainly located in highly conserved segments of the structural proteins.

3 **Authors:** Martí Cortey^{1*}, Yanli Li¹, Ivan Díaz², Hepzibar Clilverd¹, Laila Darwich^{1,2}, Enric
4 Mateu^{1,2}

5 **Affiliations:** 1 Dept Sanitat i Anatomia Animals, Facultat de Veterinària, Travessera
6 dels Turons s/n, Universitat Autònoma de Barcelona, 08193 Cerdanyola del Vallès,
7 Spain.

8 2 Centre de Recerca en Sanitat Animal (CRESA-IRTA-UAB), campus UAB,
9 Universitat Autònoma de Barcelona, 08193 Cerdanyola del Vallès, Spain.

10 *Corresponding autor

11 E-mail: marti.cortey@uab.cat

12 **Telephone:** +34 935813297

13

14 **Abstract**

15 The *Severe acute respiratory syndrome coronavirus 2* (SARS-CoV-2) pandemic offers a
16 unique opportunity to study the introduction and evolution of a pathogen into a
17 completely naïve human population. We identified and analysed the amino acid
18 mutations that gained prominence worldwide in the early months of the pandemic.
19 Eight mutations have been identified along the viral genome, mostly located in
20 conserved segments of the structural proteins and showing low variability among
21 coronavirus, which indicated that they might have a functional impact. At the moment
22 of writing this paper, these mutations present a varied success in the SARS-CoV-2 virus
23 population; ranging from a change in the spike protein that becomes absolutely
24 prevalent, two mutations in the nucleocapsid protein showing frequencies around 25%,
25 to a mutation in the matrix protein that nearly fades out after reaching a frequency of
26 20%.

27 Keywords: SARS-CoV-2, pandemia, mutation, coronavirus, fitness

28 **1. Introduction**

29 The emergence of the novel *Severe acute respiratory syndrome coronavirus 2* (SARS-
30 CoV-2) and the subsequent pandemic has become a health problem unparalleled in the
31 last century. SARS-CoV-2 is thought to be originated from an animal coronavirus that
32 successfully adapted to humans. The species of origin of SARS-CoV-2 has not been fully
33 identified, but the virus seems to be related to SARS-CoV and other coronaviruses found
34 in bats and other mammal species, although different from them (Chan et al. 2020; Lu
35 et al. 2020; Zhou et al. 2020).

36 The SARS-CoV-2 genome size is around 30 kb with the typical gene structure known in
37 other betacoronaviruses: starting from the 5', more than two-thirds of the genome
38 comprises orf1ab encoding polyproteins (nsp1 to nsp15), while the last third consists of
39 genes encoding major structural proteins; including spike (S or ORF2), envelope (E or
40 ORF4), membrane (M or ORF5), and nucleocapsid (N or ORF9) proteins. Additionally, the
41 SARS-CoV-2 contains at least 6 minor structural proteins, encoded by ORF3a, ORF6,
42 ORF7a, ORF7b, ORF8, and ORF10 genes (Khailany et al. 2020).

43 The first cases of the novel coronavirus associated disease (CoVID-19) have been traced
44 to the Chinese province of Hubei in early December 2019
45 (<https://www.who.int/csr/don/12-january-2020-novel-coronavirus-china/en/>).

46 Although the actual index case is not really known, the first sequence of the novel
47 coronavirus was produced within weeks from the emergence of the disease (Zhu et al.
48 2019). As of the moment of writing this paper, more than 16,000 sequences have been
49 produced in less than five months since the start of the pandemic. This is a unique
50 opportunity to gain insight on the evolution of a betacoronavirus in a completely naïve

51 human population. In this context, viral variants efficiently transmitted will have less
52 influence of the selection exerted by the immune response, since most transmissions
53 will occur from individuals before the development of an efficient immune response to
54 naïve recipients.

55 The aim of the present study was to determine the amino acid substitutions in viral
56 proteins that were widely present in available sequences of SARS-CoV-2, relating them
57 to the known chronology of the pandemic. Also, the mutations found were assessed in
58 order to try to understand its potential significance for viral fitness.

59 **2. Material and methods**

60 **2.1. Sequences**

61 SARS-CoV-2 sequences were retrieved from GISAID database (<https://www.gisaid.org/>).
62 The full set of sequences used in the present study included the 12,562 high-quality
63 complete sequences available on May 3rd, 2020. Additionally, a reference sequence from
64 SARS-CoV, pangolin, civet, and three from bat coronaviruses (Genbank accession
65 numbers AY278741, MT084071, AY572034, KY417146, MN996532 and MK211376,
66 respectively) were used for comparative purposes. The set of sequences were arranged
67 chronologically by date of isolation after the first reported SARS-CoV-2 sequences
68 (identified as Wuhan-01 from December 24th, 2019).

69 **2.2. Analysis of non-synonymous mutations and selection of mutations to be studied**

70 Complete genomes were aligned using the multiple alignment program ClustalW
71 (Thompson et al. 1994) and consequently split by week according to their isolation date
72 with the sequence alignment editor Bioedit (Hall 1999). Using Wuhan-01 as the
73 reference, an arbitrary date for the 1st report of the amino acid changes at the end of

74 February 2020 was set to represent an early date of the pandemic, three months after
75 the initial case was reported. Also, an arbitrary cut-off frequency of 10% was set to select
76 the amino acid substitutions that were considered widespread. Thus, any substitution
77 reported before the end of February and present in 10% or more of the frequencies in a
78 given week was studied. In order to check if the variants identified presented a
79 worldwide distribution, their geographical distribution was summarized by continents.
80 For each substitution, an alignment with the homologous proteins of SARS-CoV, civet,
81 pangolin, and bat coronaviruses was performed to assess whether the mutation
82 affected conserved or variable regions.

83 **2.3 Comparison with predominant amino acid substitutions in early, mid and late** 84 **cases of SARS epidemic of 2003-04**

85 To compare predominant non-synonymous mutations occurring during different phases
86 of the 2003 SARS-CoV epidemic, all sequences of SARS-CoV available at Genbank for
87 which date of the case was available (directly or through literature search) were
88 collected. Sequences were classified as early, mid or late based on the common
89 classifications of cases (He et al. 2004). Analysis of mutations was done similarly to SARS-
90 CoV-2.

91 **2.4 Analysis of the potential biological significance of the observed substitutions**

92 Modelling of the original S protein in Wuhan-01 and the mutant protein was produced
93 using SWISS-MODEL protein template 6vsb.1.A (<https://swissmodel.expasy.org/>).
94 Accuracy of the models was assessed by the Global Model Quality Estimation (GQME)
95 and the Qualitative Model Energy Analysis (QMEAN) scores. 3D structures were
96 rendered using PyMOL (The PyMOL Molecular Graphics System, Version 2.3.4.

97 Schrödinger, LLC). The same program was used to determine changes in the protein
98 structure or distances between atoms or residues. The set of proteomic utilities in
99 EXPASY (<https://www.expasy.org/proteomics>) was used to check for different aspects
100 on the mutant proteins (motifs, phosphorylation sites, etc.). PROVEAN 1.1.
101 (<http://provean.jcvi.org/>, Choi et al. 2012) was used to gain insight on whether the
102 mutation could be deleterious or neutral. Changes in the secondary structure of proteins
103 were predicted by using the CFSSP: Chou & Fasman Secondary Structure Prediction
104 Server (<https://www.biogem.org/tool/chou-fasman/>, Ashok 2013).

105 When mutations affected known epitopes, impact on antigenicity was evaluated by
106 means of epitope prediction tools in the IBDB Resource web (<http://tools.iedb.org/>).

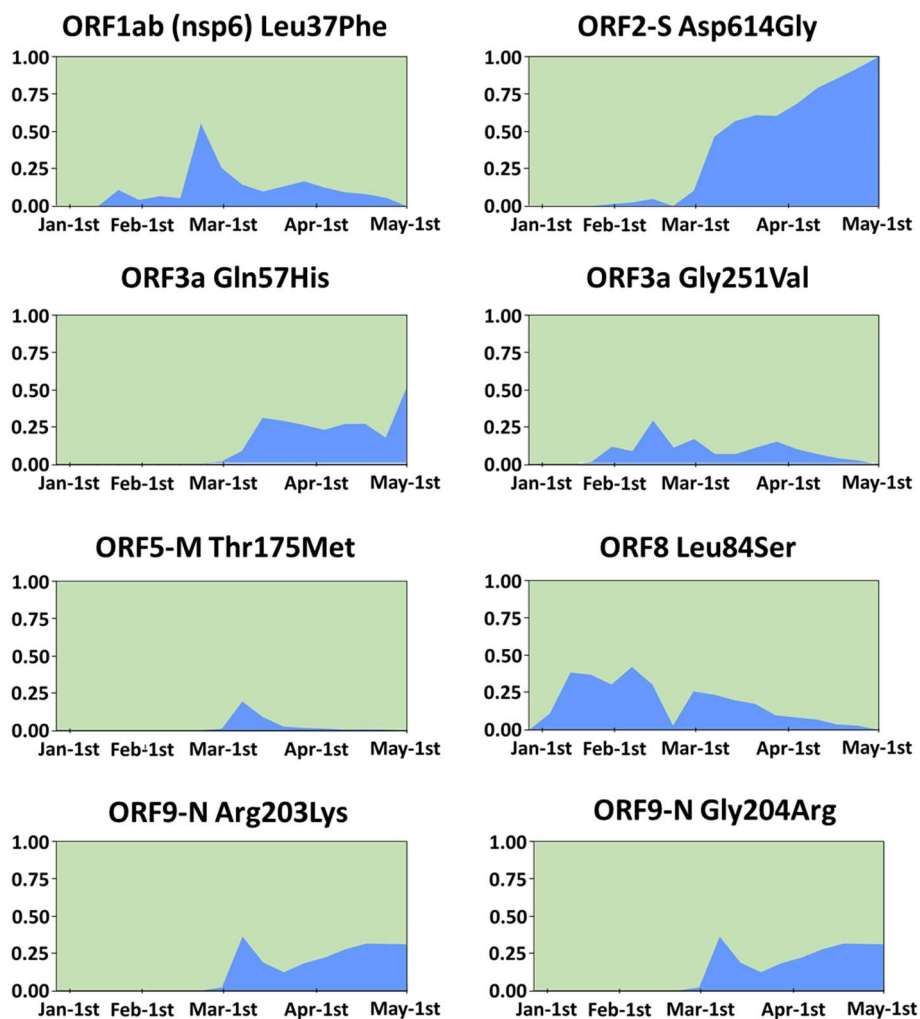
107 **3 Results**

108 **3.1. Amino acid substitutions with significant spread in the population were mainly** 109 **located in conserved segments of structural proteins of SARS-CoV-2**

110 The analysis of the set of sequences revealed that only 8 amino acid substitutions across
111 the viral genome appeared before the end of February 2020 and gained prevalence over
112 10% of the known isolates at a given time point, measuring time in weeks after the first
113 available sequence (Fig. 1). Of these, 7 substitutions were in the structural proteins
114 (ORF2-S, ORF3a, ORF5-M, ORF8, and ORF9-N) and one in the ORF1ab, specifically in
115 nsp6. Concerning their location and date, four appeared in China during January. The
116 other four appeared in Europe during the second fortnight of February, but within a
117 week, they were also reported in other continents (Supplementary Table S1).
118 Interestingly, the only substitution that became fully predominant was Asp614Gly in the
119 spike protein (ORF2-S). Gly57His in ORF3a reached a frequency of 50% at the moment

120 of writing this paper. It is worth noting that when the sequences were analysed by
121 continents, all mutations were spread worldwide, except the 175Met in the ORF5-M,
122 that was absent in Africa (Supplementary Table S1).

123 **Figure 1.** Temporal trends in the emergence and prevalence of the amino acid
124 substitutions that appeared before the end of February 2020 and were present in >10%
125 of SARS-CoV-2 sequences at any time between December 2019 and May 2020. The X-axis
126 represents time measured as weeks since December 24th, 2019 and the Y-axis
127 represents the proportion of known sequences harbouring a given mutation.



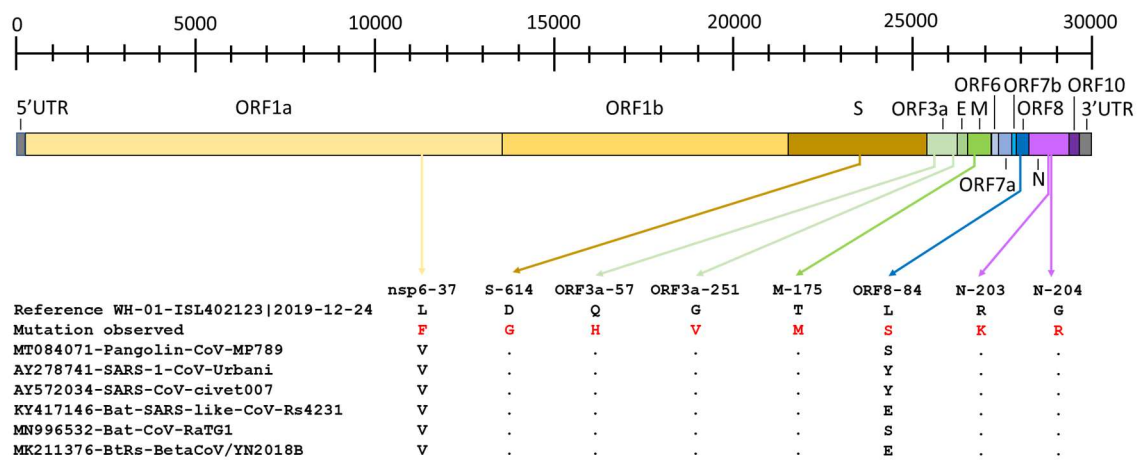
128

129 Next, we examined whether these substitutions were located in variable or conserved
130 regions of the viral genome. Interestingly, most of them corresponded to residues that

131 were conserved in SARS and related betacoronaviruses of pangolins, civets, or bats (Fig.
132 2).

133 The comparison with non-synonymous mutations that gained predominance in SARS-
134 CoV showed a different pattern. From early to late phases of the SARS epidemic of 2003,
135 11 substitutions gained wide spread. Three of them were located in nsp3, four in nsp4,
136 one in nsp16, and three in the spike protein (Supplementary Table S2). However,
137 mutations in SARS-CoV were located in conserved positions of civet and bat-related
138 coronaviruses but, those positions were different in pangolin.

139 **Figure 2.** Location of the mutations found in the present study and corresponding amino
140 acids in pangolin, SARS-CoV, civet, and bat-related coronaviruses. One-letter code is
141 used to represent amino acids. A dot is used to indicate a conserved residue compared
142 to the first sequence in the alignment (SARS-CoV-2 isolate Wuhan-01).



143

144

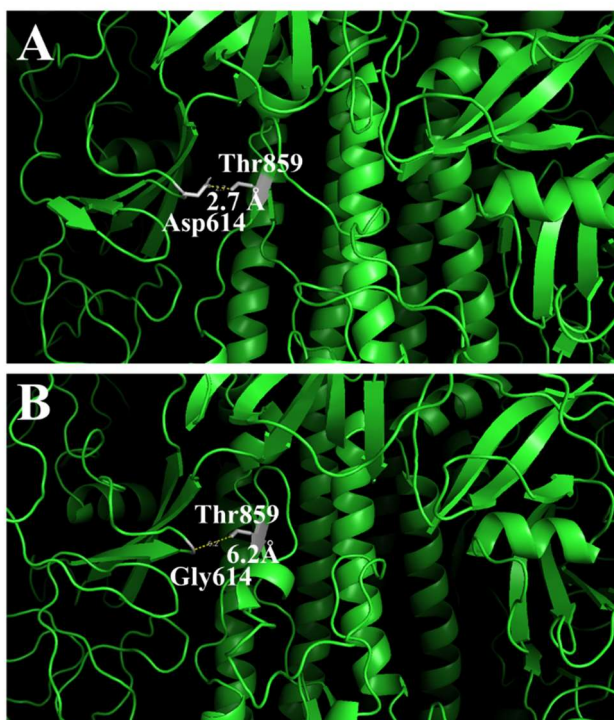
145 3.2. Substitutions in the major structural proteins S, N and M

146 The examination of the spike protein sequences of SARS-CoV-2 revealed that the
147 Asp614Gly mutation that appeared in January 2020 in Shanghai, gained predominance
148 with time. Thus, by the end of April 2020 it was almost 100% prevalent (Fig. 1).

149 Next, we checked whether this mutation emerged in a particular SARS-CoV-2 clade or
150 its spread was irrespective of the clade where a given isolate could be allocated. The
151 Gly614 could be found in different branches across the phylogenetic tree of SARS-CoV-
152 2 (available at <https://nextstrain.org/>).

153 Residue 614 is located in the S1 domain of the spike protein. Modelling of the original
154 and mutant proteins using template 6vsb.1.A from SwissModel of the S protein of SARS-
155 CoV-2 (Wrap et al. 2020) showed that by changing Asp614 by Gly614, the distance to
156 Thr859 and its side chain increased (from 2.7 to 6.2 Å), creating a cavity and a more
157 relaxed structure (Fig. 3). The analysis of the antigenicity of the potential epitope
158 homologous to 597-625 of SARS-CoV did not show any significant difference between
159 mutants.

160 **Figure 3.** Changes in the 3-D structure of the S-protein in the original (A) and mutant (B)
161 proteins. The pictures show amino acid residues on 20Å distance from Asp614 (A) or
162 Gly614 (B) and their distance to Thr859.



163

164 In the nucleocapsid phosphoprotein, a double mutation (Arg203Lys-Gly204Arg) was
165 observed to gain predominance during the pandemic (Fig. 1). This mutation is in the
166 serin-rich (SR) segment of the linker region (LKR) of the protein. The N protein of SARS-
167 CoV is known to be bound by Ubc9, a ubiquitin conjugating enzyme of the sumoylation
168 system, probably at the SR segment. Since lysines are targets for sumoylation, we next
169 checked whether this substitution could introduce a sumoylation motif. The prediction
170 using two different methods JASSA v4 (<http://www.jassa.fr/index.php?m=jassa>) and
171 GPS-SUMO (<http://sumosp.biocuckoo.org/>) was negative. Predicted phosphorylation
172 sites and enzymes (NetPhos 3.1.) were not significantly different between variants. The
173 observed mutations were predicted to be neutral by Provean.

174 In the matrix protein, the main change was the substitution of Thr175 by Met175.
175 Thr175 was predicted to form part of a motif known to interact with 14-3-3 proteins and
176 of a potential phosphorylation site (173-SRTLSYYKL-181) targeted by protein kinases A
177 (PKA) and C (PKC), ribosomal s6 kinase (RKS), and DNA-dependent protein kinases. The
178 substitution of the THR by a MET was predicted again to be a phosphorylation site (for
179 PKA, PKC and RKS). Provean predicted that the introduction of M175 was deleterious
180 (score -3.135). Since no reliable model was available, 3D structure could not be assessed.
181 Interestingly, the presence of a M175 in the M protein occurred together with the
182 203KR204 mutation in 98% of the cases ($p < 0.0001$). The rapid decay of this mutation
183 (Fig. 1) would be consistent with a deleterious effect.

184 **3.3. Substitutions in the minor structural proteins 3a and 8, and the non-structural**
185 **protein nsp6**

186 Two substitutions were identified in the protein encoded by ORF3a. The first was a
187 substitution of Gln57 present in the Wuhan-01 isolate by His57, and the second was the
188 substitution of Gly251 by Val251. To note, only 0.04% of the sequences harboured the
189 double His57/Val251 mutation. The presence of His in residue 57 would be expected to
190 result in an increased positive charge at that site.

191 The Gly251Val mutation occurred in a predicted serine-phosphorylation site 248-
192 TID(**G/V**)SSGVV-256. The introduction of a Val reduced the prediction scores for the site
193 from a maximum >0.90 with Gly for phosphokinase B and ATM serine/threonine
194 phosphokinase to 0.73-0.82 for the same enzymes. It is worth to note that the amino
195 acid at this position was strongly correlated with the amino acid present in position 57
196 of the same protein. Thus, among the sequences harbouring His57 in ORF3a, 99.8% were
197 associated with Gly251 and only 0.2% with Val251. In contrast, for Gln57, 17% of the
198 sequences harboured Val251 ($p < 0.001$). Similarly, His57 was only found in sequences
199 harbouring Arg203-Gly204 in the N protein, while Gln57 was found simultaneously with
200 Arg203-Gly204 (70% of the cases) or Lys203-Arg204 (30%) ($p < 0.001$). Both mutations
201 were predicted as deleterious (scores of -3.286 for His57 and -8.581 for Val251).

202 The lack of a reliable model made impossible to make any prediction of the impact of
203 those mutations on the 3D structure or the interactions between residues. However,
204 the checking of the potential in the secondary structure of the protein revealed that
205 mutation Gln57His resulted in the elimination of a turn of the protein predicted to be in
206 Ser58. Similarly, mutation Gly251Val eliminated the turn at Gly251 but did not affect the
207 turn predicted for Ser253 (Supplementary Fig. S3).

208 Regarding the ORF8 protein, it is worth noting that the mutation of residue 84
209 (Leu84Ser) happened simultaneously with a silent mutation in nucleotide position 8987
210 (ORF1ab, nsp4, A→T). This permitted to distinguish 2 clades in the initial weeks of the
211 pandemic that contained isolates from Wuhan, Shanghai, and Hong-Kong
212 (Supplementary Fig. S4). These clades did correspond to the L and S types reported by
213 Tang et al. (2020). This mutation was predicted to be neutral.

214 Finally, the last change was found in the nsp6 protein, Leu37Phe, which significance was
215 unclear. This mutation was also predicted to be neutral.

216 **4. Discussion**

217 The present SARS-CoV-2 pandemic is a worst-case scenario of the introduction of a new
218 agent that transmits easily in a completely naïve population. In this context,
219 transmission events occur in an uncommonly high scale, in a very short period of time
220 and with little selective pressures from the immune system if compared to an endemic
221 situation. This scenario would permit the arising of a great diversity of viral variants of
222 which the fitter could be expected to gain predominance.

223 In the present study, we identified 8 early mutations in the SARS-CoV-2 genome that
224 gained prevalence over 10% at some point during the pandemic. This cut-off was
225 arbitrarily set to discriminate random mutations and errors in sequencing from changes
226 that might have a bigger impact. Certainly, this approach has the limitation of neglecting
227 some mutations with lesser prevalence that still can be biologically significant. Time will
228 show it.

229 It is worth noting that 7 out of 8 of the widely spread mutations occurred in residues
230 that were highly conserved in related coronaviruses of bats, pangolins, civets, or in SARS-

231 CoV (Fig. 2). Conserved regions are usually assumed to be functionally relevant and thus,
232 mutations in them may have deleterious effects or can be hardly tolerated; if so, they
233 will be probably removed in the future. A mutation in a highly conserved region that
234 becomes widespread and persists can be thought as representative of a change that
235 increases viral fitness. In the present case, we found three different situations:
236 mutations that expanded and rise to predominance, mutations that expanded to a
237 certain extent and fade out, and mutations that are apparently expanding but not yet
238 predominant. This pattern affecting conserved regions was also seen for SARS-CoV
239 although the affected proteins were different. Interestingly, in SARS-CoV-2 most of the
240 mutations were in structural proteins, while in SARS-CoV were in non-structural ones,
241 suggesting that the adaption process from the original host species to human was
242 different in these two cases. When the spike protein was examined, this difference was
243 more obvious. Mutations in SARS-CoV occurred in positions conserved in the civet and
244 bat-related coronaviruses but different from those of pangolin and SARS-CoV-2. In
245 contrast, spike mutation in position 614 of SARS-CoV-2 affected a residue that was
246 conserved in betacoronaviruses of pangolins, civets, bats, and SARS-CoV. This would be
247 compatible with a scenario where those mutations affected the viral fitness for that
248 particular new host, namely humans. The scale of the viral replication in the scenario of
249 a pandemic would be unprecedented for those coronaviruses and will provide the
250 probability for such beneficial mutations to appear and expand.

251 The Asp614Gly mutation in S protein is an example of a mutation becoming fully
252 predominant. A previous report (Korber et al. 2020) already indicated its emergence.
253 Recently, Bhattacharyya et al. (2020) suggested that the predominance of this mutation
254 as the pandemic advanced and the low proportion in initial phases of it was related to a

255 single nucleotide deletion in the transmembrane protease serine 2 (TMPRSS2) that is
256 common in Europeans and North Americans but rare in East Asians. The Asp614Gly
257 mutation would introduce a cleavage site for that enzyme. This would explain, at least
258 partially, the spread of this mutant outside Asia.

259 We do agree with Korber et al. (2020) regarding the possibility that the Asp214Gly
260 mutation produced a laxer interaction between S1 and S2 spike domains that might
261 facilitate shedding of S1 from membrane-bound S2. In SARS-CoV, the segment of S
262 protein including residues 597-625 contains epitopes inducing both neutralizing
263 antibodies (positions 604-625 in SARS-CoV) and antibodies participating in an antibody-
264 dependent enhancement (ADE) in animal models (residues 597-625) (Wang et al. 2016).
265 The mutation Asp614Gly would affect the epitope segment inducing ADE but not the
266 one inducing neutralizing antibodies. Although, it has been hypothesized that ADE may
267 have a role in CoVID-19 (Tetro 2020) this has not been demonstrated (Sharma 2020).
268 The effects of the mutation on the immune escape or the transmission potential cannot
269 be concluded at this moment.

270 Regarding the mutations found in the nucleocapsid phosphoprotein, the first surprising
271 fact was to find two consecutive substitutions in the highly conserved serine-rich
272 segment of the LKR region of the protein. The LKR region is essential for conferring
273 flexibility to the protein as well as for cell signalling, binding to RNA and to M protein. It
274 also contains multiple phosphorylation sites in the SR segment that are thought to be
275 essential (Reviewed by McBride et al. 2014). Since the mutant nucleocapsid introduced
276 a Lys residue - a canonical target for sumoylation - and sumoylation occurs in this region
277 by analogy to SARS-CoV (Fan et al. 2006) we also tested it. No significant differences

278 were determined between the original and the mutant sequences using prediction tools
279 for phosphorylation or sumoylation.

280 Interaction of the nucleocapsid protein with the M protein is thought to happen
281 between the SR-region and the C-terminal domain of M (Escors et al. 2001; Kuo and
282 Masters 2002). We found that the M175 phenotype of the M protein was almost
283 exclusively associated to the 203K-204R phenotype of the nucleocapsid protein. It is
284 tempting to hypothesize that the above-mentioned residues may be involved in such
285 interaction. Besides this, the introduction of an additional charge in the SR segment may
286 enhance interaction with RNA which core is negatively charged.

287 Changes in the ORF3a protein have been recently reported to define microclonal clades
288 of SARS-CoV2 (Issa et al. 2020). We have found that the introduction of Val251 was
289 apparently non compatible with His57. We determined that, probably, such mutations
290 affected the secondary structure by changing the number of turns. Considering that
291 both residues are out of the functional domains proposed by Issa et al. (2020), those
292 changes in the structure of the protein may modify interactions enough to be non-
293 compatible. Interestingly, mutations in ORF3a protein and mutations in the
294 nucleocapsid protein were related. To our knowledge, these two proteins have not been
295 investigated for interactions; however, this finding suggests that might be an interaction
296 between them.

297 ORF8 mutation Leu84Ser was reported before (Tang et al. 2020). The authors suggested
298 that the Leu variant (called L) is more aggressive and spreads easier. The evolution of
299 the proportion of strains harbouring this mutation would not support the idea of a
300 higher transmissibility of the L variant since its frequency clearly declined. Certainly, the

301 introduction of such variant in different countries applying control measures earlier or
302 later could have had an impact on the spread of the variants as well.

303 Finally, the nsp6 mutation at position 37 is difficult to interpret. The nsp6 of
304 coronaviruses is part of the replication machinery of the virus and has been reported to
305 induce autophagosomes (Benvenuto et al. 2020). The same authors suggested that it
306 may lead to a lower stability of the protein structure. Lacking a verified model for the
307 protein it is difficult to assess whether this happens or not. Interestingly, the same
308 authors indicated that the distribution of the Leu phenotype was restricted to Asia while
309 the Phe was common in other parts of the world. According to our analysis, the Phe
310 phenotype has almost disappeared in current sequences. The discrepancy could be
311 originated in the fact that Benvenuto et al. (2020) analysed the 351 sequences available
312 in a past moment while in the present analysis thousands of sequences from all over the
313 world have been included.

314 In summary, the present is a comprehensive report of the amino acid mutations that
315 gained spread during the SARS-CoV-2 pandemic up to now. Most of the substitutions
316 gaining wide diffusion occurred in conserved positions indicating that they probably had
317 a functional impact but, differently from SARS-CoV, they accumulated in structural
318 proteins. Interestingly, most of these mutations faded out, except for the Asp614Gly in
319 the S protein that became predominant suggesting that it contributed to viral fitness.
320 Some others are still increasing in prevalence, like the mutations in the nucleocapsid
321 protein here reported and that might be related to the mutations in the M protein. This
322 data may serve to gain further insight in the evolution of SARS-CoV-2.

323 5. Acknowledgements

324 The authors would like to thank to all those who contribute in the fight against the SARS-
325 CoV-2. We would like to thank as well all those who kindly shared genome data in
326 publicly available databases, available for download from <https://www.gisaid.org/>.

327 6. Author contributions

328 MC, YL, ID, HC, LD and EM participated in all tasks of the present work.

329 7. Funding

330 Martí Cortey was funded by the Ramon y Cajal program (reference RYC-2015-17154).
331 Hepzibar Clilverd was supported by a fellowship of the Spanish Ministry of Science and
332 Innovation (program FPU). No other funding sources.

333 **8. Conflict of interest:** None declared.

334 9. Supplementary Material

335 **Supplementary Table S1.** Distribution by continents, in percentages, of the 8 amino acid
336 mutations detected across the 12,562 genomes analysed.

337

338 **Supplementary Table S2.** Localizations of the amino acid mutations in SARS CoV that
339 gained prevalence compared with SARS-CoV-2, civet, pangolin, and bat-related
340 coronaviruses.

341

342 **Supplementary Figure S3.** Secondary structure of the original ORF3a (A) or the mutant
343 variants (B & C).

344

345 **Supplementary Figure S4.** Neighbor-Joining tree based on the Tamura-Nei distances,
346 constructed with the 20 earliest complete genomes harbouring a Leucine or a Serine in
347 the amino acid residue 84 of the ORF8 protein. Isolation dates as reported in the GISAID
348 database (www.gisaid.org). Numbers along the internal branches represent their
349 confidence after the initial dataset was resampled with 1,000 bootstrap replicates.

350 **10. References**

351 Ashok-Kumar, T. (2013) 'CFSSP: Chou and Fasman Secondary Structure Prediction
352 server', *Wide Spectrum*, 9: 15-19.

353 Benvenuto, D. et al. (2020) 'Evolutionary analysis of SARS-CoV-2: how mutation of Non-
354 Structural Protein 6 (NSP6) could affect viral autophagy', *Journal of Infection*, 4453:
355 30186–9.

356 Bhattacharyya , C. et al. (2020) 'Global Spread of SARS-CoV-2 Subtype with Spike Protein
357 Mutation D614G is Shaped by Human Genomic Variations that Regulate Expression of
358 TMPRSS2 and MX1 Genes', *bioRxiv*, doi.org/10.1101/2020.05.04.075911.

359 Chan J.F., et al. (2020) 'Genomic characterization of the 2019 novel human-pathogenic
360 coronavirus isolated from a patient with atypical pneumonia after visiting Wuhan',
361 *Emerging Microbes & Infections*, 9: 221–36.

362 Choi, Y. et al. (2012) 'Predicting the Functional Effect of Amino Acid Substitutions and
363 Indels', *PLoS ONE*, 7: e46688.

- 364 Escors D. et al. (2001) 'The membrane M protein carboxy terminus binds to transmissible
365 gastroenteritis coronavirus core and contributes to core stability', *Journal of Virology*,
366 75: 1312–24.
- 367 Fan, Z. et al. (2006) 'SARS-CoV nucleocapsid protein binds to hUbc9, a ubiquitin
368 conjugating enzyme of the sumoylation system', *Journal of Medical Virology*, 78: 1365-
369 73.
- 370 Hall, T. A. (1999) 'BioEdit: a user-friendly biological sequence alignment editor and
371 analysis program for Windows 95/98/NT', *Nucleic Acids Symposium Series*, 41: 95–8.
- 372 He, J. F. et al. (2004) 'Molecular Evolution of the SARS Coronavirus During the Course of
373 the SARS Epidemic in China', *Science*, 303: 1666-9.
- 374 Issa, E. et al. (2020) 'SARS-CoV-2 and ORF3a: Non-Synonymous Mutations and
375 Polyproline Regions', *mSystems*, 5: e00266–20.
- 376 Khailany, R. A. et al. (2020) 'Genomic characterization of a novel SARS-CoV-2', *Gene*
377 *Reports*, 19: 100682.
- 378 Korber B. et al. (2020) 'Spike mutation pipeline reveals the emergence of a more
379 transmissible form of SARS-CoV-2', *bioRxiv* 2020.04.29.069054.
- 380 Kuo, L. and Masters, P.S. (2002) 'Genetic evidence for a structural interaction between
381 the carboxy termini of the membrane and nucleocapsid proteins of mouse hepatitis
382 virus', *Journal of Virology*, 76: 4987–99.
- 383 Lu, R. et al. (2020) 'Genomic characterisation and epidemiology of 2019 novel
384 coronavirus: implications for virus origins and receptor binding', *Lancet*, 395: 565–74.

- 385 McBride, R. et al. (2014) 'The coronavirus nucleocapsid is a multifunctional protein',
386 *Viruses*, 6: 2991–3018.
- 387 Sharma, A. (2020) 'It is too soon to attribute ADE to COVID-19', *Microbes and Infection*,
388 S1286-4579(20)30051-4.
- 389 Tang, X. et al. (2020) 'On the origin and continuing evolution of SARS-CoV-2', *National*
390 *Science Review*, nwaa036.
- 391 Tang, X. et al. (2020) 'On the origin and continuing evolution of SARS-CoV-2', *National*
392 *Science Review*, nwaa036.
- 393 Tetro, J.A. (2020) 'Is COVID-19 receiving ADE from other coronaviruses?', *Microbes &*
394 *Infection*, 22: 72–3.
- 395 Thompson, J. D. et al. (1994) 'CLUSTAL W: improving the sensitivity of progressive
396 multiple sequence alignment through sequence weighting, position-specific gap
397 penalties and weight matrix choice', *Nucleic Acids Research*, 22: 4673–80.
- 398 Wang, Q. et al. (2016) 'Immunodominant SARS Coronavirus Epitopes in Humans Elicited
399 both Enhancing and Neutralizing Effects on Infection in Non-human Primates', *ACS*
400 *Infectious Diseases*, 5: 361–76.
- 401 Wrapp, D. et al. (2020) 'Cryo-EM structure of the 2019-nCoV spike in the prefusion
402 conformation', *Science*, 367: 1260–3.
- 403 Zhou, P. et al. (2020) 'A pneumonia outbreak associated with a new coronavirus of
404 probable bat origin', *Nature*, 579: 270–3.

405 Zhu, N. et al. (2020) 'A Novel Coronavirus from Patients with Pneumonia in China, 2019',
406 *The New England Journal of Medicine*, 382: 727–33.

407

408 **Supplementary Table S1.** Distribution by continents, in percentages, of the 8 amino acid mutations detected across the 12,562 SARS-CoV-2
 409 genomes analysed.

Location	nsp6-37		S-614		ORF3a-57		ORF3a-251		M-175		ORF8-84		N-203		N204	
1st Report	Wuhan 18/01		Shanghai 23/01		France 21/02		Hong-Kong 23/01		Netherlands 24/02		Wuhan 05/01		England 23/02		England 23/02	
Aa Change	Leu	Phe	Asp	Gly	Gln	His	Gly	Val	Thr	Met	Leu	Ser	Arg	Lys	Gly	Arg
Asia (n=741)	81,8	18,2	69,1	30,9	82,8	17,2	90,6	9,4	98,6	1,4	80,6	19,4	94,4	5,6	94,4	5,6
Europe (n=7,571)	81,1	18,9	20,6	79,4	83,3	16,7	86,1	13,9	92,0	8,0	93,2	6,8	68,6	31,4	68,6	31,4
America (n=3,238)	95,8	4,2	37,2	62,8	48,6	51,4	97,9	2,1	99,4	0,6	72,0	28,0	96,1	3,9	96,1	3,9
Africa (n=110)	90,3	9,7	14,7	85,3	80,0	20,0	96,8	3,2	100,0	0,0	92,6	7,4	90,4	9,6	90,4	9,6
Oceania (n=902)	69,7	30,3	46,5	53,5	69,9	30,1	87,4	12,6	95,5	4,5	79,5	20,5	86,7	13,3	86,7	13,3
Mean (n=12,562)	84,2	15,8	29,7	70,3	73,1	26,9	89,7	10,3	94,6	5,4	86,1	13,9	78,7	21,3	78,7	21,3

410

411

412 **Supplementary Table S2.** Localizations of the amino acid mutations in SARS CoV that gained prevalence compared with SARS-CoV-2, civet,
 413 pangolin, and bat-related coronaviruses.

Protein	ORF1ab (nsp3)	ORF1ab (nsp3)	ORF1ab (nsp3)	ORF1ab (nsp4)	ORF1ab (nsp4)	ORF1ab (nsp4)	ORF1ab (nsp4)	ORF1ab (nsp16)	ORF2S	ORF2S	ORF2S
Aa position in the protein	303	844	1298	6	231	307	332	138	242	347	1166
Original SARS-CoV	T	I	F	W	A	A	A	K	L	R	E
Mutated Aa in SARS-CoV	I	L	L	C	V	V	V	R	S	K	K
AY572034-SARS-CoV-civet007	T	I	L	W	A	A	A	R	S	R	E
SARS-CoV-2 WH-01-ISL402123	T	I	A	W	A	A	V	K	T	R	K
MT084071-Pangolin-CoV-MP789	T	I	Del	W	A	A	V	K	T	R	K
KY417146-Bat-SARS-like-CoV-Rs4231	T	I	L	W	A	A	A	K	L	R	K
MN996532-Bat-CoV-RaTG13	T	I	A	W	A	A	A	K	T	R	K
MK211376-BtRs-BetaCoV/YN2018B	T	I	L	W	A	A	A	K	P	R	K

414

415 **Supplementary Figure S3. Secondary structure of the original ORF3a (A) or the mutant**
416 **variants (B & C). H = alpha-helix, E = Beta sheet, T = Turn**

417 **A) Original ORF3a protein**
418
419
420
421
422
423
424
425
426

		10	20	30	40	50	60		
		*	*	*	*	*	*		
Query	1	MDLFMRIFTIGTVTLKQGEIKDATPSDFVRATATIPIQASLPPGWLIVGVALLAVFQ	SASKIITLKKRWQ	70					
Helix	1	HHHHHHHHHHHHHHHHHHHHHH	HH	70					
Sheet	1	EEEEEEEEEEEEEEEE	EEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE	70					
Turns	1		TT T T	T			T	70	
Struc	1	HHHEEEEEEEEEEH	HHHHHHHT	CEEEEEHHEE	HHHT	EEEEEEEEEE	HHEEHHT	HHHEEHHHHEHH	70

427
428
429
430
431
432
433
434
435

		220	230	240	250	260	270		
		*	*	*	*	*	*		
Query	211	YYQLYSTQLSTDTGVEHVTF	FFIYNKIVDEPEEHVQIHTIDGSSG	VVNPVMEPIYDEPTTT	TSVPL	275			
Helix	211	HHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHH	HHHHHHHH	275					
Sheet	211	EEEEEEEEEEEEEEEEEEEEEEEEEEEE	EEEE	EEEEEEEEEEEEEEEEEEEE	275				
Turns	211	T	T T	T T		T	275		
Struc	211	EEEEEEEEEEEEEH	HEEEEEEE	HHHHHHHEE	EEHCTCT	CEEEHHHEE	HEET	EEEEEECCC	275

436
437
438 **B) Gln57His mutant**
439
440
441
442
443
444
445
446
447

		10	20	30	40	50	60		
		*	*	*	*	*	*		
Query	1	MDLFMRIFTIGTVTLKQGEIKDATPSDFVRATATIPIQASLPPGWLIVGVALLAVFQ	SASKIITLKKRWQ	70					
Helix	1	HHHHHHHHHHHHHHHHHHHHHH	HH	70					
Sheet	1	EEEEEEEEEEEEEEEE	EEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE	70					
Turns	1		TT T T	T			T	70	
Struc	1	HHHEEEEEEEEEEH	HHHHHHHT	CEEEEEHHEE	HHHT	EEEEEEEEEE	HHEEH	HHHHHEEHHHHEHH	70

448
449
450 **C) Gly251Val mutant**
451
452
453
454
455
456
457
458
459

		220	230	240	250	260	270		
		*	*	*	*	*	*		
Query	211	YYQLYSTQLSTDTGVEHVTF	FFIYNKIVDEPEEHVQIHTIDGSSG	VVNPVMEPIYDEPTTT	TSVPL	275			
Helix	211	HHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHH	HHHHHHHH	275					
Sheet	211	EEEEEEEEEEEEEEEEEEEEEEEEEEEE	EEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE	275					
Turns	211	T	T T	T		T	275		
Struc	211	EEEEEEEEEEEEEH	HEEEEEEE	HHHHHHHEE	EEHCTCT	CEEEHHHEE	HEET	EEEEEECCC	275

460
461
462

463 **Supplementary Figure S4.** Neighbor-Joining tree based on the Tamura-Nei distances,
464 constructed with the 20 earliest complete genomes harbouring a Leucine or a Serine in
465 the amino acid residue 84 of the ORF8 protein. Isolation dates as reported in the GISAID
466 database (<https://www.gisaid.org/>). Numbers along the internal branches represent
467 their confidence after the initial dataset was resampled with 1,000 bootstrap replicates.

