

# 1 SARS-CoV-2 Point Mutation and Deletion Spectra, and Their

## 2 Association with Different Disease Outcome

3 Brenda Martínez-González<sup>1</sup>, María Eugenia Soria<sup>1,2</sup>, Lucía Vázquez-Sirvent<sup>1</sup>, Cristina  
4 Ferrer-Orta<sup>3</sup>, Rebeca Lobo-Vega<sup>1</sup>, Pablo Mínguez<sup>4,5,6</sup>, Lorena de la Fuente<sup>4,5,6</sup>, Carlos  
5 Llorens<sup>7</sup>, Beatriz Soriano<sup>7</sup>, Ricardo Ramos<sup>8</sup>, Marta Cortón<sup>4,5</sup>, Rosario López-  
6 Rodríguez<sup>4,5</sup>, Carlos García-Crespo<sup>2</sup>, Isabel Gallego<sup>2,9</sup>, Ana Isabel de Ávila<sup>2</sup>, Jordi  
7 Gómez<sup>9,10</sup>, Luis Enjuanes<sup>11</sup>, Llanos Salar-Vidal<sup>1</sup>, Jaime Esteban<sup>1</sup>, Ricardo Fernandez-  
8 Roblas<sup>1</sup>, Ignacio Gadea<sup>1</sup>, Carmen Ayuso<sup>4,5</sup>, Javier Ruíz-Hornillos<sup>12,13,14</sup>, Nuria  
9 Verdaguer<sup>3</sup>, Esteban Domingo<sup>2,9\*</sup> and Celia Perales<sup>1,9,11\*</sup>

10 <sup>1</sup>Department of Clinical Microbiology, Instituto de Investigación Sanitaria-Fundación  
11 Jiménez Díaz University Hospital, Universidad Autónoma de Madrid (IIS-FJD, UAM)  
12 Av. Reyes Católicos 2, 28040 Madrid, Spain, <sup>2</sup>Centro de Biología Molecular “Severo  
13 Ochoa” (CSIC-UAM), Consejo Superior de Investigaciones Científicas (CSIC),  
14 Campus de Cantoblanco, 28049 Madrid, Spain, <sup>3</sup>Structural Biology Department,  
15 Institut de Biologia Molecular de Barcelona CSIC, 08028 Barcelona, Spain,  
16 <sup>4</sup>Department of Genetics & Genomics, Instituto de Investigación Sanitaria-Fundación  
17 Jiménez Díaz University Hospital, Universidad Autónoma de Madrid (IIS-FJD, UAM),  
18 Av. Reyes Católicos 2, 28040 Madrid, Spain, <sup>5</sup>Centre for Biomedical Network Research  
19 on Rare Diseases (CIBERER), Instituto de Salud Carlos III, 28029 Madrid, Spain,  
20 <sup>6</sup>Bioinformatics Unit, Instituto de Investigación Sanitaria-Fundación Jiménez Díaz  
21 University Hospital, Universidad Autónoma de Madrid (IIS-FJD, UAM), Madrid  
22 28040, Spain, <sup>7</sup>Biotechvana, “Scientific Park”, Universidad de Valencia, 46980  
23 Valencia, Spain, <sup>8</sup>Unidad de Genómica, “Scientific Park of Madrid”, Campus de  
24 Cantoblanco, 28049 Madrid, Spain, <sup>9</sup>Centro de Investigación Biomédica en Red de  
25 Enfermedades Hepáticas y Digestivas (CIBERehd), Instituto de Salud Carlos III, 28029  
26 Madrid, Spain, <sup>10</sup>Instituto de Parasitología y Biomedicina ‘López-Neyra’ (CSIC),  
27 Parque Tecnológico Ciencias de la Salud, Armilla, 18016 Granada, Spain,  
28 <sup>11</sup>Department of Molecular and Cell Biology, Centro Nacional de Biotecnología (CNB-  
29 CSIC), Consejo Superior de Investigaciones Científicas (CSIC), Campus de  
30 Cantoblanco, 28049 Madrid, Spain, <sup>12</sup>Allergy Unit, Hospital Infanta Elena, Valdemoro,  
31 Madrid, Spain; <sup>13</sup>Instituto de Investigación Sanitaria-Fundación Jiménez Díaz  
32 University Hospital, Universidad Autónoma de Madrid (IIS-FJD, UAM), Av. Reyes  
33 Católicos 2, 28040 Madrid, Spain; <sup>14</sup>Faculty of Medicine, Universidad Francisco de  
34 Vitoria, Madrid, Spain

35  
36 \*Corresponding authors: Celia Perales ([cperales@cbm.csic.es](mailto:cperales@cbm.csic.es)) and Esteban Domingo  
37 ([edomingo@cbm.csic.es](mailto:edomingo@cbm.csic.es))

38 **ABSTRACT**

39 Mutant spectra of RNA viruses are important to understand viral pathogenesis, and  
40 response to selective pressures. There is a need to characterize the complexity of mutant  
41 spectra in coronaviruses sampled from infected patients. In particular, the possible  
42 relationship between SARS-CoV-2 mutant spectrum complexity and disease  
43 associations has not been established. In the present study, we report an ultra-deep  
44 sequencing (UDS) analysis of the mutant spectrum of amplicons from the nsp12  
45 (polymerase)- and spike (S)-coding regions of thirty nasopharyngeal isolates (diagnostic  
46 samples) of SARS-CoV-2 of the first COVID-19 pandemic wave (Madrid, Spain, April  
47 2020) classified according to the severity of ensuing COVID-19. Low frequency  
48 mutations and deletions, counted relative to the consensus sequence of the  
49 corresponding isolate, were overwhelmingly abundant. We show that the average  
50 number of different point mutations, mutations per haplotype and several diversity  
51 indices was significantly higher in SARS-CoV-2 isolated from patients who developed  
52 mild disease than in those associated with moderate or severe disease (exitus). No such  
53 bias was observed with RNA deletions. Location of amino acid substitutions in the three  
54 dimensional structures of nsp12 (polymerase) and S suggest significant structural or  
55 functional effects. Thus, patients who develop mild symptoms may be a richer source of  
56 genetic variants of SARS-CoV-2 than patients with moderate or severe COVID-19.

57

58 **IMPORTANCE**

59 The study shows that mutant spectra of SARS-CoV-2 from diagnostic samples differ in  
60 point mutation abundance and complexity, and that significantly larger values were  
61 observed in virus from patients who developed mild COVID-19 symptoms. Mutant  
62 spectrum complexity is not a uniform trait among isolates. The nature and location of  
63 low frequency amino acid substitutions present in mutant spectra anticipate great  
64 potential for phenotypic diversification of SARS-CoV-2.

65 **Keywords:** COVID-19 severity, Mutant spectrum, Diversity index, Mutation, Deletion,  
66 nsp12 (polymerase), spike, Ultra-deep sequencing.

## 67 INTRODUCTION

68 Betacoronavirus SARS-CoV-2 emerged in the human population in 2019, and it  
69 is the causal agent of the new pandemic disease COVID-19 (1), with a death toll which  
70 is increasing at the time of this writing (<https://covid19.who.int/>). Genetic variations in  
71 SARS-CoV-2 genomes [annotated in the GISAID (<https://www.gisaid.org/>), PubMed  
72 (<https://www.ncbi.nlm.nih.gov/pmc/>), and ENA data banks  
73 (<https://www.ebi.ac.uk/ena/browser/home>); among others] affect non-structural and  
74 structural protein-coding regions. Despite the short history of SARS-CoV-2 circulation,  
75 newly arising variants exhibiting different mutational patterns are regularly being  
76 identified. A distinction has been made between variants of interest (VOI), due to  
77 features with potential impact (such as transmissibility), and variants of concern (VOC),  
78 due to definite evidence of enhanced transmissibility  
79 (<https://www.who.int/en/activities/tracking-SARS-CoV-2-variants/>). New SARS-CoV-2  
80 variants are likely to become prominent as COVID-19 continues, despite natural or  
81 vaccine-induced immunity (2-5). Likewise, the generation of viral escape mutants is a  
82 major concern as a potential limitation of immune and antiviral agent efficacy for  
83 SARS-CoV-2 (6-10), as it has been established for other RNA viruses.

84 The first step in the diversification of viruses during their epidemic spread is the  
85 generation of variants within each infected host. This pattern of intra-host evolution  
86 results in the formation of mutant spectra that constitute reservoirs of genetic and  
87 phenotypic virus variants in the infected host (11, 12). Studies with several RNA viruses  
88 have shown that viral intra-mutant spectrum complexity, estimated by the average  
89 number of mutations per genome, expressed by a series of diversity indices [Shannon  
90 entropy, maximum mutation frequency, Gini Simpson, nucleotide diversity, number of  
91 polymorphic sites, and number of haplotypes (13, 14)] may have an impact on viral  
92 tropism, viral persistence, disease progression and response to antiviral interventions  
93 [several cases have been described or reviewed in (11, 15-22)]. Evidence of  
94 quasispecies dynamics has been reported for SARS-CoV-2 (23-29), as well as for other  
95 coronaviruses (30-34). However, it is unclear how mutant spectrum complexity  
96 parameters of this emerging pathogen vary among different viral isolates, and whether  
97 previously observed effects of mutant spectrum composition on RNA virus behavior  
98 apply also to SARS-CoV-2, particularly its connection with disease severity.

99 Two recent studies indicated higher mutant spectrum complexity in SARS-CoV-  
100 2 from patients who developed severe disease than mild disease, either analyzing the  
101 spike (S)-coding regions (35), or the entire genome with limited mutant spectrum  
102 resolution (36). In the present study, we have examined mutant spectra of the nsp12  
103 (polymerase)- and S-coding regions of SARS-CoV-2 present in 30 nasopharyngeal  
104 swab samples taken at the time of diagnosis of patients progressing towards disparate  
105 disease outcomes. Applying a 0.5% cut-off value for point mutation and deletion  
106 detection, using SeekDeep as bioinformatics platform, we found that virus from patients  
107 who developed mild disease exhibited a significantly higher mutant spectrum  
108 complexity than virus from patients who developed moderate or severe disease (exitus).  
109 The difference occurred both in the nsp12 (polymerase)- and S-coding regions. In  
110 contrast, no significant differences in the spectrum of minority deletions were observed  
111 among virus from the three patient's categories (mild, moderate or severe disease).  
112 Some amino acid substitutions found at low frequency in mutant spectra, including  
113 substitutions with low statistical acceptability and with potential functional effects, are  
114 nevertheless present in SARS-CoV-2 isolates recorded in data banks.

115

## 116 **RESULTS**

117 **SARS-CoV-2 mutant spectra from patients progressing towards different COVID-**  
118 **19 severity.** We previously classified 448 patients [Fundación Jiménez Díaz (FJD)  
119 cohort, Madrid, Spain, April 2020] according to the COVID-19 severity into three  
120 categories: mild, moderate and severe COVID-19 —based on a number of demographic  
121 and clinical parameters— and we found a positive association between viral load in  
122 nasopharyngeal swabs and disease severity (37). For the present study, we have chosen  
123 thirty of the nasopharyngeal samples based on three criteria: (i) the COVID-19  
124 category, including 10 patients who developed mild symptoms, 10 patients who  
125 developed moderate disease, and 10 patients who progressed to severe disease and  
126 exitus; (ii) patients whose diagnostic (RT-PCR RNA samples) displayed similar Ct  
127 values (average  $Ct=25.37\pm 3.9$  for mild,  $Ct=21.81\pm 2.4$  for moderate, and  $Ct=20.38\pm 2.9$   
128 for exitus patients); and (iii) similar time interval between symptom onset and swab  
129 collection (average  $5.78\pm 4.2$  days for mild,  $4.89\pm 3.1$  days for moderate and  $4.5\pm 2.6$   
130 days for exitus patients). When present, comorbidities were equally represented among

131 the different COVID-19 severities (Table S1 in  
132 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>).

133 To set up ultra-deep sequencing (UDS) analyses of SARS-CoV-2 obtained from  
134 nasopharyngeal swabs, we have adapted experimental protocols previously used for  
135 HCV quasispecies characterization (38-41), and applied the SeekDeep pipeline (42) to  
136 the analysis of minority point mutations and deletions in SARS-CoV-2 mutant spectra  
137 (described in Materials and Methods). RNA from nasopharyngeal swabs was extracted,  
138 amplified and subjected to UDS using MiSeq platform (Illumina). Four amplicons (A1  
139 to A4) covering nucleotides 14,534 to 16,054 of the nsp12 (polymerase)-coding region,  
140 and two amplicons (A5 and A6) covering nucleotides 22,872 to 23,645 of the S-coding  
141 region were analyzed (Fig. 1). The total number of clean reads was 19,592,197,  
142 corresponding to 653,073 (range 316,710-910,727) reads per patient, that yielded an  
143 average of 110,689 (range 38,865-215,662) clean reads per amplicon, with a 0.5% cut-  
144 off frequency for point mutations and deletions (Fig. S1 in  
145 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>).

146 To provide a general picture of SARS-CoV-2 divergence and mutant spectrum  
147 heterogeneity, we constructed a heat map representing the frequency of each variation  
148 in the nsp12 (polymerase) and S-coding regions (point mutations and deletions; no  
149 insertions were detected), relative to the genomic sequence of a Wuhan isolate  
150 (identified as NCBI reference sequence NC\_045512.2), and divided the samples  
151 according to different COVID-19 severity (Fig. 2 and Table S2 in  
152 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>). Considering all patients analyzed,  
153 the number of positions that included a variation (either a point mutation or a deletion)  
154 was two-fold higher in the S-coding region (105 positions with a genomic modification  
155 out of 774 positions analyzed) than in the nsp12 (polymerase)-coding region (91  
156 positions modified out of 1,521 positions analyzed). In addition to minority mutations in  
157 each mutant spectrum, a total of six different dominant mutations relative to the  
158 reference sequence (those with frequencies between 90% and 100%) were also present;  
159 they are identified as “Divergence” in Fig. 2. This class of mutations has been excluded  
160 for the quantification of mutations and complexity indices in a mutant spectrum.  
161 Ninety-four percent of mutations were found at frequencies that ranged between 0.5%  
162 and 30% within its mutant spectrum, whereas only 6% corresponded to “Divergence”  
163 mutations ( $p < 0.001$ ; proportion test). Interestingly, 62 out of 97 point mutations (64%)  
164 within the mutant spectra were detected at frequencies below 2% (Fig. 2).

165 To evaluate if some parameters of the mutant spectra (considering only point  
166 mutations present at a frequency below 30%) were associated with COVID-19 severity,  
167 we first counted the number of different point mutations present in virus from each  
168 patient group. In the two coding regions analyzed, the average number of different  
169 mutations in virus from patients with mild disease was significantly higher than in virus  
170 from patients with moderate disease or exitus [ $p < 0.001$  for the comparison between  
171 mild versus moderate and mild versus exitus, both for nsp12 (polymerase)- and S-  
172 coding region; proportion test]; no significant difference was noted between moderate  
173 and exitus patients [ $p = 0.081$  and  $p = 0.603$  for nsp12 (polymerase)- and S-coding  
174 regions, respectively; proportion test]; normalization of the number of different  
175 mutations to the length of the regions analyzed did not modify the result (Fig. 3A). No  
176 such difference among patient groups was observed with the number of different  
177 deletions (all  $p$ -values  $> 0.05$ ; proportion test), although a trend towards a larger number  
178 of deletions in virus from patients who developed mild disease was maintained in the S-  
179 coding region (Fig. 3B). Thus, SARS-CoV-2 mutant spectra from diagnostic samples of  
180 patients who evolved to mild disease included a significantly larger average number of  
181 mutations, but not of deletions, than virus from patients who progressed towards  
182 moderate or severe (exitus) COVID-19.

183

184 **Evaluation of complexity indices.** The comparison of SARS-CoV-2 mutant spectra  
185 was extended to two groups of diversity indices: abundance (which consider the reads  
186 of entities and their frequency in the mutant spectrum), and incidence (which consider  
187 only reads of entities) (13). To this aim, we have adapted the QSutils package (43) to  
188 the quantification of diversity indices for SARS-CoV-2 mutant spectra (described in  
189 Materials and Methods). In the nsp12 (polymerase)-coding region, a significant increase  
190 of the values of abundance and incidence indices was observed in samples from patients  
191 who developed mild disease, as compared with samples from patients with moderate  
192 disease ( $p < 0.001$  for  $H_S$ ,  $H_{GS}$ ,  $Mf_{max}$  and  $\pi$ ;  $p=0.001$  for number of polymorphic sites  
193 and number of haplotypes; Wilcoxon test). Also significant was the difference between  
194 samples associated with mild disease and severe disease (exitus) ( $p = 0.004$  for  $H_S$ ,  $p =$   
195  $0.010$  for  $H_{GS}$ ,  $p = 0.012$  for  $Mf_{max}$  and  $p = 0.010$  for  $\pi$ ;  $p = 0.004$  for number of  
196 polymorphic sites and number of haplotypes; Wilcoxon test). The same tendency was  
197 observed in the S-coding region but the differences did not reach statistical significance  
198 (all  $p$ -values  $> 0.05$ ; proportion test) (Fig. 4 and Table S3 in

199 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>). In each amplicon, a larger number  
200 of haplotypes was found in samples associated with mild than moderate or severe  
201 disease, and the majority of mutated haplotypes included only one mutation (Fig. S2 in  
202 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>). Thus, the higher abundance of  
203 mutations in SARS-CoV-2 mutant spectra from patients who exhibited only mild  
204 symptoms is also reflected in an increase of mutant spectrum complexity.

205

### 206 **Point mutation and amino acid substitution types in SARS-CoV-2 mutant spectra.**

207 Considering mutant spectra of all samples analyzed, transitions and non-synonymous  
208 mutations were more abundant than transversions and synonymous mutations,  
209 respectively, with different degrees of statistical significance (Table 1); a similar trend  
210 was also observed when the samples were divided according to COVID-19 severity of  
211 the patients.

212 In the nsp12 (polymerase)-coding region, the frequency of mutation types  
213 normalized to base composition ranked as follows: T to C > A to G > C to T; when  
214 dividing the samples according to disease severity, the most frequent mutation in exitus  
215 patients was C to T (Fig. S3A in <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>). In  
216 the S-coding region the ranking was T to C > A to G = C to T (Fig. S3B in  
217 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>). T to C transitions were the most  
218 frequent mutation type in the third codon position (67.50%), whereas A to G was the  
219 most prevalent type at the second and first codon positions (45.16% and 38.46%,  
220 respectively).

221 The amino acid substitutions found in nsp12 (polymerase) and S were  
222 positioned in the three-dimensional structure of the proteins [Protein Data Bank  
223 (<http://www.rcsb.org/>)], their statistical acceptability was evaluated with PAM250  
224 matrix (44), and their potential functional effects was estimated by applying the SNAP2  
225 predictor (45). All amino acid substitutions found in nsp12 (polymerase) and S are  
226 listed in Table S2 in <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>, together with  
227 their PAM250 and SNAP2 scores; their location in the three dimensional structure of  
228 the proteins is depicted in Fig. S4 in <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>.  
229 Those amino acid substitutions which suggest alteration of protein structure or function  
230 are described in Tables 2 and 3. Some of the substitutions in nsp12 (polymerase) predict  
231 positive or negative functional effects (Table 2 and Fig. 5). For example, V557I may

232 enhance the stability of the interaction with nitrogen base T+1, and Q822H predicts  
233 increased stability of loop in the thumb domain. In contrast, D618N abolishes the  
234 catalytic aspartate of polymerase in domain A, and C765R should distort the catalytic  
235 domain (Table 2 and Fig. 5). The amino acid substitutions observed in S tend to  
236 increase the hydrophobicity of the region where they are located (Table 3 and Fig. 6).  
237 The replacement of A by V at position 475 may enhance interactions of S with ACE2;  
238 A522V may contribute to stabilize the RBD domain in the “open” position through  
239 contacts with neighbor V, T, P and L residues; R567G could facilitate fusion with the  
240 host cell; A570V may bring closer two S chains (Table 3 and Fig. 6). Drastic  
241 substitutions may belong to defective genomes that have a transient existence or that  
242 may be maintained by complementation (see Discussion).

243

244 **Deletion repertoire in SARS-CoV-2 mutant spectra.** Deletions were also analyzed by  
245 UDS with a cut-off value of 0.5% (as detailed in Materials and Methods), with the same  
246 reads used for point mutations. The analyses identified five different deletions which  
247 spanned 3-13 nucleotides (nt) in the nsp12 (polymerase)-coding region, and five  
248 different deletions that spanned 2-51 nt in the S-coding region (Figs. 2 and S5 in  
249 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>). In the nsp12 (polymerase)-coding  
250 region, the 4 nt and 13 nt deletions that disrupted the coding frame generated a stop  
251 codon 10 and 26 residues downstream, respectively. The 2 nt, 16 nt, 22 nt, and 28 nt  
252 deletions in the S-coding region led to stop codons 3 to 18 nucleotides downstream (Fig.  
253 S5 in <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>). The number of deletions that  
254 generated a stop codon was significantly higher in the S-coding region (26 out of 27  
255 deletions) than in nsp12 (polymerase)-coding region (2 out of 10 deletions) ( $p < 0.001$ ;  
256 proportion test). The sites of deletions did not map in homopolymeric regions or tandem  
257 repeats, and they were not flanked by the same nucleotide types (Fig. S5 in  
258 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>).

259

260 **Point mutation and deletion hot spots.** The distribution of genomic variations (point  
261 mutations and deletions) per amplicon was similar for the four amplicons of the nsp12  
262 (polymerase)-coding region ( $p$ -value  $> 0.05$ ; proportion test). In contrast, amplicon A6  
263 of the S-coding region accumulated higher number of total mutations than A5 ( $p < 0.001$ ;  
264 proportion test) (Fig. 7A). This difference may result from dissimilar functional  
265 constraints on the protein portions represented by each amplicon, i.e. a uniform



266 distribution of polymerase motifs A to G among the four nsp12 (polymerase)  
267 amplicons, compared with the presence of the receptor-binding domain (RBD) in  
268 amplicon A5 of S (compare Figs. 1 and 7A).

269 Hot spots for SARS-CoV-2 variations have been described based on the  
270 comparison of consensus sequences of independent isolates (46-48). Here we have  
271 defined as hot spots those positions that presented the same point mutation or deletion  
272 in the mutant spectrum of at least five different isolates (Fig. 7B and Table S2 in  
273 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>). Two hot spots were located in the  
274 nsp12 (polymerase)-coding region (a point mutation at position 15,756, and a deletion  
275 of residues 14,856 to 14,858), and two in the S-coding region (a point mutation at  
276 position 23,544 and a deletion of residues 23,555 to 23,582) (Fig. 7B). These hot spots  
277 do not coincide with those reported for SARS-CoV-2 consensus sequences (46-48).

278

279 **Geographical and temporal characterization of mutations based on CoV-GLUE**  
280 **database.** SARS-CoV-2 mutant spectra from infected patients can include mutations  
281 that are also found as dominant in later isolates (27). In the mutant spectra of the 30  
282 samples from our cohort, the ratio of amino acid substitutions (including those  
283 corresponding to divergence mutations) that were unique [not yet annotated in the CoV-  
284 GLUE database that is enabled by GISAID metadata (49)] versus those described in  
285 other (prior or subsequent) isolates was 0.2 (10 out of 60). Out of the 60 non-  
286 synonymous mutations, 8 (13.33%) were described worldwide at about the same time  
287 that they were identified in our cohort, and 19 (31.67%) were described afterwards (Fig.  
288 S6 in <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>). Of particular interest is S  
289 protein substitution S494P, located at the ACE-2 binding region (Table S2 in  
290 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>) that reached epidemiological  
291 importance, and was found in some isolates of the alpha variant. Thus, SARS-CoV-2  
292 mutant spectra—in particular from patients that developed mild symptoms— may  
293 constitute a rich reservoir of mutations with the potential to be represented in  
294 epidemiologically relevant variants.

295

## 296 **DISCUSSION**

297 The UDS analysis of the nsp12 (polymerase)- and S-coding regions of 30  
298 biological samples without cell culture passage confirm the presence of complex SARS-

299 CoV-2 mutant spectra in diagnostic nasopharyngeal samples of the virus (23-28).  
300 Contrary to a previous conclusion with other patient cohorts (35, 36), our  
301 quantifications show that in both the nsp12 (polymerase)- and the S-coding regions  
302 analyzed there was a positive association between the number of point mutations and a  
303 mild disease manifestation in the corresponding patients. No such association was  
304 observed with the minority deletions that also populated the mutant spectra (Figs. 2 and  
305 3). There are several non-mutually exclusive mechanisms that may contribute to a larger  
306 average number of point mutations in samples from patients that developed mild disease  
307 than in those from patients with moderate or severe disease. One is that the major sites  
308 of replication of the virus may not be identical in the three groups of patients.  
309 Mutational input may be affected by a variety of host cell functions, including editing  
310 activities (50), or as a consequence of the effects on polymerase fidelity of non-  
311 structural viral proteins that participate in genome replication, as evidenced with other  
312 RNA viruses (51-54). This possibility for SARS-CoV-2 is suggested by non-identical  
313 preferred transition mutation types in the isolates, depending on the associated disease  
314 severity (Table 1). A second influence may lie in a longer time of asymptomatic intra-  
315 host virus replication prior to the onset of mild symptoms and COVID-19 diagnosis. A  
316 prolonged replication time does not necessarily imply a larger viral load in the infected  
317 host, but it may entail an increase in the average number of variant genomes in the  
318 population. Another possibility is that bottleneck events—which may transiently reduce  
319 the number of mutations scored within mutant spectra—intervene with higher intensity  
320 in patients doomed to severe disease than those developing mild disease. This may  
321 come about through the immune response that may partially suppress viral replication,  
322 and that it is also part of the COVID-19 pathogenesis process (55-57). Several  
323 possibilities may explain dissimilar conclusions with other studies: (i) independent  
324 cohorts may have been infected by virus belonging to clades displaying non-identical  
325 behavior, and (ii) methodological differences such as in the criteria to classify patients  
326 according to COVID-19 symptoms, in the PCR-UDS resolution attained, or in the  
327 sample type taken for analysis (naso/oropharyngeal swabs versus nasopharyngeal  
328 aspirates), among others. The multiple factors that contribute to a mutant spectrum  
329 complexity beg for studies with other cohorts to try to clarify whether complexity of  
330 viral RNA in diagnostic samples responds to discernible virological parameters, and  
331 whether UDS data might help predicting disease evolution or response to treatment, as  
332 previously documented for hepatitis C (58, 59).

333 We have focused the mutant spectrum analysis on two regions of the SARS-  
334 CoV-2 genome whose encoded proteins are likely subjected to widely different  
335 constraints. The nsp12 (polymerase) is involved in genome replication and  
336 transcription, and the S glycoprotein has a major role in virus attachment, fusion and  
337 entry, as well as in defining the antigenic profile of the virus. A total of 41 different  
338 amino acid substitutions in nsp12 and 15 substitutions in S have been recorded in the 30  
339 mutant spectra analyzed (Table S2 in  
340 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>). Normalization to the sequenced  
341 protein length gives an average frequency of non-synonymous mutations of 8% for  
342 nsp12 and 6% for S in the mutant spectra. Three substitutions in S map in the receptor  
343 binding domain (RBD). One of them, A475V (present at 26% frequency in virus from a  
344 patient who developed mild disease) reduced the sensitivity to several monoclonal  
345 antibodies (60). S494P (dominant in virus from a patient who developed mild disease)  
346 was listed among the nine most frequent substitutions in a large-scale study of 506,768  
347 SARS-CoV-2 isolates; it is considered a likely vaccine-escape substitution, and possibly  
348 involved also in increased transmissibility of some isolates of the alpha variant detected  
349 beginning September 2020 (61-63) (<https://www.cdc.gov/>).

350 Substitutions that are present at low frequency are associated with predicted  
351 more drastic structural and functional effects, and, some of them have been identified in  
352 the sequences compiled in the CoV-GLUE data base (compare Table 2 and Fig. S6 in  
353 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>).

354 It is likely that disruptive amino acid substitutions belong to defective or  
355 minimally replicating (very low fitness) genomes that have either a transient existence  
356 in the population or that can be maintained at detectable levels by complementation (64)  
357 (for example those with lesions incompatible with polymerization activity). Defective  
358 genomes need not represent a biological or evolutionary dead end. They can exert  
359 modulatory effects on the entire population (65), and they also constitute a rich  
360 substrate for RNA recombination to rescue viable genomes that may become  
361 epidemiologically competent viruses.

362 Newly replicated genomes *in vivo* may incorporate deletions as a result of  
363 limited processivity of the coronavirus replicase (66, 67). Genomes with deletions may,  
364 on average, be subjected to stronger negative selection than genomes with point  
365 mutations, blurring differences in their frequency among samples from the three patient  
366 categories. This is likely to apply mainly to out of frame deletions that give rise to

367 truncated proteins; for example, in the S-coding region we have identified deletion  
368 ( $\Delta$ )23,555 to 23,570,  $\Delta$ 23,555 to 23,582, and  $\Delta$ 23,561 to 23,582 which are located near  
369 the S1/S2 cleavage site, and are expected to impair S function. Their maintenance to the  
370 point of reaching sufficient concentration to be detectable by UDS may reflect a higher  
371 efficiency of complementation of *trans*-acting structural proteins than non-structural  
372 proteins (64). This may also explain the lower frequency of out of frame deletions in the  
373 nsp12 (polymerase)- than in the S-coding region. It has been proposed that defective S  
374 proteins generated around the S1/S2 cleavage site could potentially reduce the severity  
375 of the infection (68).

376 All point mutations and deletions were found at frequencies below 30% in the  
377 corresponding mutant spectra. Several important biological and clinical features could  
378 influence the shape of SARS-CoV-2 mutant spectra. However, it should be considered  
379 that the large size of the coronavirus genome may limit the accumulation of mutations  
380 relative to less complex RNA genomes, due to negative effects of mutations on fitness  
381 (69). Not even the point mutation hot spots were found at frequencies above 1% in the  
382 quasispecies where they were present (compare Figs. 2 and 7). This is compatible with  
383 hot spots reflecting sites where lesions are more tolerated within a generally constrained  
384 RNA genome. The fact that hot spots according to mutant spectra do not coincide with  
385 those defined by consensus sequences adds to other observations that indicate that  
386 residue conservation criteria at these two levels do not coincide (70). That the great  
387 majority of mutations in SARS-CoV-2 mutant spectra are present at low frequency may  
388 slow down the response of the virus to specific selective constraints such as inhibitors  
389 or neutralizing antibodies. Under this scenario, viral load may become more important  
390 to furnish genomes with mutations required to respond to the constraints (71).  
391 Comparative measurements with different RNA viruses are needed to endorse these  
392 potential effects of mutant spectrum composition.

393 The higher percentage of transitions versus transversions, and of non-  
394 synonymous versus synonymous mutations is in agreement with previous reports of  
395 mutant spectrum and consensus sequence analyses of SARS-CoV-2 (23, 24, 35, 68, 72).  
396 Some differences with previous studies have been observed in the preferred mutation  
397 types (Fig. S3 in <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>). While in the  
398 mutant spectra of our cohort T to C was the most frequent point mutation, other studies  
399 reported C to T as the preferred mutation type (72, 73). C to T was, however, the most  
400 frequent mutation in virus from the subset of exitus patients (Fig. S3 in

401 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>) hinting at the possibility that in  
402 previous studies virus from patients with moderate and severe COVID-19 might have  
403 been over-represented. The lack of dominance of C to U transitions in our samples is  
404 also reflected in absence of depletion of amino acids A, H, Q, P and T when considering  
405 all amino acid substitutions observed (50, 74); the data of Fig. S3 in  
406 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5> show a net gain of 3 amino acids in  
407 the A, H, Q, P and T subset. Another possible explanation for differences with previous  
408 studies could be that the latter focused on consensus sequences obtained from data bases  
409 covering the whole genome, whereas our results correspond to two specific genomic  
410 regions sequenced by UDS.

411 The six point mutations that altered the consensus sequence of the mutant  
412 spectra relative to reference NC\_045512.2 (identified as “Divergence” in the heat map  
413 of Fig. 2 and in Table S2 in <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>) allowed  
414 an estimate of the rate of accumulation of mutations in the SARS-CoV-2 consensus  
415 sequence. The time interval between our Madrid isolates (dated April 2020) and the  
416 reference Wuhan isolate (dated December 2019) was 4 months. Considering this time  
417 interval, the average rate of evolution calculated is  $(1.6 \pm 0.6) \times 10^{-3}$  mutations per  
418 nucleotide and year (m/n/y), and it is only slightly higher than the average value from  
419 ten previous studies  $(1.2 \pm 0.6) \times 10^{-3}$  m/n/y (range  $9.9 \times 10^{-4}$  to  $2.2 \times 10^{-3}$  m/n/y) (73,  
420 75-83). Higher evolutionary rates are frequently obtained the shorter is the time interval  
421 between the virus isolations considered for the calculation [reviewed in (84)]. The  
422 values for SARS-CoV-2 are comparable to those reported for other RNA viruses,  
423 suggesting that constraints at the quasispecies level may not affect significantly  
424 evolutionary rates considered at the epidemiological level (85). Our results hint at the  
425 possibility that SARS-CoV-2 evolving in patients exhibiting mild symptoms may  
426 contribute a majority of the variants that drive the high rates of evolution quantified at  
427 the epidemiological level.

428

## 429 MATERIALS AND METHODS

430 **Patient cohort and stratification.** Samples were collected during the first COVID-19  
431 outbreak in Spain. The cohort of the study included 30 patients admitted to the  
432 Fundación Jiménez Díaz Hospital (FJD, Madrid, Spain) from April 3 to 29, 2020. All  
433 patients were confirmed to be positive for SARS-CoV-2 by a specific real-time RT-

434 PCR (VIASURE Real Time PCR) with a Ct (cycle threshold, which is inversely  
435 correlated with viral RNA level) range of 15.6 to 28.5; the samples are a subset from the  
436 cohort that has been previously described in (37). Data collected included patient  
437 demographics, risk factors for COVID-19, and clinical information at the time of  
438 SARS-CoV-2 diagnosis (Table S1 in  
439 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>). The parameters used to classify the  
440 patients included: (i) need of hospitalization, (ii) need of mechanical ventilation, (iii)  
441 admission to the intensive care unit (ICU), and (iv) exitus attributed to COVID-19.  
442 Taking these parameters into account, the patients were classified as mild, moderate and  
443 severe (exitus) cases according to the symptoms and hospitalization requirements: (i)  
444 mild symptoms (neither hospital admission nor ICU) (n=10), (ii) moderate symptoms  
445 (hospitalization without ICU) (n=10), and (iii) severe symptoms (hospitalization with  
446 admission to the ICU, and progression to exitus in all cases) (n=10). The clinical  
447 classification was established before the data analysis was performed.

448 **Oligonucleotide design.** To design oligonucleotide primers, a total of 663 SARS-CoV-  
449 2 sequences from the NCBI database ([https://www.ncbi.nlm.nih.gov/genbank/sars-cov-](https://www.ncbi.nlm.nih.gov/genbank/sars-cov-2-seqs/)  
450 [2-seqs/](https://www.ncbi.nlm.nih.gov/genbank/sars-cov-2-seqs/)) were retrieved and aligned to the Wuhan-Hu-1 NCBI reference sequence  
451 NC\_045512.2 (86). Nucleotide sequences were analyzed to design forward and reverse  
452 oligonucleotide primers (Table S4 in <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>).  
453 Four pairs of oligonucleotides were used for amplification and sequencing of four  
454 overlapping amplicons of the genomic region of nsp12 (polymerase) (nucleotides  
455 14,511 to 16,075) encoding amino acids 366 to 871, and two pairs to cover the region of  
456 the S protein (nucleotides 22,853 to 23,666) encoding amino acids 438 to 694 (residue  
457 numbering according to reference sequence NC\_045512.2) (Fig. 1 and Table S5 in  
458 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>).

459

460 **RNA extraction and amplification of SARS-CoV-2 RNA from infected patients.**

461 SARS-CoV-2 RNA was extracted from 140 µl of medium from nasopharyngeal swabs  
462 using the QIAamp Viral RNA Mini Kit (250) (Qiagen), as specified by the  
463 manufacturer. Amplifications of nsp12 (polymerase)- and S-coding regions were  
464 performed by RT-PCR. Each region was amplified from 5 µl of the RNA preparation by  
465 RT-PCR using Transcriptor One Step RT-PCR kit (Roche Applied Science). To  
466 perform the RT-PCR, 5 µl of the preparation were mixed with 10 µl of 5x buffer, and 2

467  $\mu\text{l}$  of a solution containing the forward primer, 2  $\mu\text{l}$  of a solution with the reverse primer  
468 (50 ng/ $\mu\text{l}$ , each), and 1  $\mu\text{l}$  of polymerase. Reaction parameters were 50°C for 30 min for  
469 the reverse transcription, an initial denaturing step at 94°C for 7 min, followed by 35  
470 cycles of a denaturing step at 94°C for 10 s, an annealing step at 46-48°C for 30 s, an  
471 extension step at 68°C for 40 s, and then a final extension at 68°C for 7 min. In the case  
472 of samples with a Ct value greater than 26 (6 samples from the mild symptom group),  
473 the number of cycles was increased to 45. Negative controls (amplification reactions in  
474 the absence of RNA) were included in parallel to ascertain absence of contamination by  
475 template nucleic acids. Amplification products were analyzed by 2% agarose gel  
476 electrophoresis, using Gene Ruler 1 Kb Plus DNA ladder (Thermo Scientific) as molar  
477 mass standard. PCR products were purified (QIAquick Gel Extraction Kit, Qiagen),  
478 quantified (Qubit dsDNA Assay kit, ThermoFisher Scientific), and tested for quality  
479 (TapeStation System, Agilent Technologies) prior to sequencing using the Illumina  
480 MiSeq platform. Dilutions of 1:10, 1:100 and 1:1,000 of the initial RNA preparation  
481 and subsequent amplification by RT-PCR were carried out for one patient of each  
482 disease severity (Fig. S7 in <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>). When  
483 amplification with the 1:1,000 dilution of template produced a visible DNA band, the  
484 ultra-deep sequencing analysis was performed with the undiluted template to avoid  
485 redundant copying of the same template molecules, as we have previously documented  
486 (87, 88).

487

488 **Ultra-Deep Sequencing of SARS-CoV-2 from Infected Patients.** PCR products were  
489 adjusted to  $4 \times 10^9$  molecules/ $\mu\text{l}$  before generating DNA pools that were purified using  
490 Kapa Pure Beads (KapaBiosystems, Roche), and quantified using Qubit as previously  
491 described (38-40), and then fixed at 1.5 ng/ $\mu\text{l}$ . Purified DNA pools were further  
492 processed using the DNA library preparation kit Kapa Hyper Prep kit (Roche), during  
493 which each pool was indexed using SeqCap Adapter Kit A/B (Nimblegen) (24 Index).  
494 Each DNA pool was quantified by LightCycler 480, and sequenced using MiSeq  
495 sequencing platform with MiSeq Reagent kit v3 ( $2 \times 300$  bp mode with the 600 cycle  
496 kit) (Illumina).

497

498 **Bioinformatics analyses.** Controls to establish the basal error, the frequency of PCR-  
499 induced recombination, and the similarity of the results with different amplifications  
500 and sequencing runs were previously performed (38, 41, 89). Therefore, mutations

501 identified with a frequency above the 0.5% cut-off value and with coverage greater than  
502 10,000 reads were considered for the analyses, based on different controls carried out  
503 with hepatitis C virus (HCV), as detailed elsewhere (38, 90).

504 Beginning with the Fastq data, two bioinformatic pipelines [SeekDeep (42), and  
505 a new previously described pipeline for HCV (38)] were applied to HCV (Fig. S8 in  
506 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>), and then adapted to SARS-CoV-2 to  
507 quantify deletions (termed VQS-Haplotyper, freely available in Github at this address  
508 <https://github.com/biotechvana/VQS-haplotyper>) (Fig. S9 in  
509 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>). As control with an independent set  
510 of UDS data, we compared the point mutations and their frequencies within HCV  
511 quasispecies obtained using both bioinformatics procedures, and the results were very  
512 similar ( $r=0.9957$  and  $p<0.0001$ ; Pearson correlation test) (Fig. S8 in  
513 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>). For SARS-CoV-2 mutant spectra,  
514 the analysis of clean reads using both pipelines yielded a robust similar number of point  
515 mutations and their frequencies ( $r=1$  and  $p<0.0001$ ; Pearson correlation test). Also, both  
516 pipelines produced similar results for deletions and their frequencies ( $r=0.4932$  and  
517  $p=0.0011$ ; Pearson correlation test) (Fig. S9 in  
518 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>). SeekDeep was applied using the  
519 following options: `--extraExtractorCmds=-- checkRevComplementForPrimers --`  
520 `primerNumOfMismatches 3" "--extraProcessClusterCmds=--fracCutOff 0.005 --`  
521 `rescueExcludedOneOffLowFreqHaplotypes"` (42). In the present study, point mutations,  
522 deletions and their frequencies were reported using SeekDeep, and diversity indices  
523 were calculated using VQS-Haplotyper followed by QSutils (43).

524

525 **Statistics.** The correlation between results obtained by the bioinformatics pipelines was  
526 calculated using Pearson's correlation. The statistical significance of difference between  
527 the number and type of mutations in mild, moderate and exitus patients as well as the  
528 differences between type of nucleotide changes and between PAM250 (accepted point  
529 mutations 250) and SNAP2 (Screening for Non-Acceptable Polymorphisms 2) values  
530 for amino acid substitutions were calculated by the proportion test. Statistics were  
531 inferred using software R version 4.0.2. The normality of data was tested with the  
532 Shapiro-Wilk normality test and the statistically significance of differences between  
533 diversity indices was calculated with a Wilcoxon test using GraphPad Prism 8.00.

534



535 **Data availability.** The reference accession numbers of sequences retrieved from NCBI  
536 used to design oligonucleotide primers are given in Table S4 in  
537 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>. Fastq files of SARS-CoV-2 samples  
538 included in the patient cohort are available in ENA under project id “PRJEB48766”.  
539 Nucleotide and amino acid replacements in SARS-CoV-2 from infected patients have  
540 been compiled in Table S2 in <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>.

541

542 **Ethics approval and consent to participate.** This study was approved by the Ethics  
543 Committee and the Institutional Review Board of the FJD hospital (no. PIC-087-20-  
544 FJD).

545

## 546 **ACKNOWLEDGEMENTS**

547 We acknowledge all personnel in the Clinical Microbiology Department of the  
548 FJD for help with the sample and data collection. We thank all health-care professionals  
549 who attended to COVID-19 patients, and collected the clinical samples that were  
550 included in this study in a difficult moment of the COVID-19 epidemic in Spain. We  
551 thank José María Aguado and Octavio Carretero for their support to the whole project.  
552 We are indebted to Cristina Villaverde for her technical expertise and help with the  
553 samples. We acknowledge Dres. J. Gregori and J. Quer for their contribution to the  
554 quasispecies analyses of HCV-infected samples.

555 This work was supported by Instituto de Salud Carlos III, Spanish Ministry of  
556 Science and Innovation (COVID-19 Research Call COV20/00181), and co-financed by  
557 European Development Regional Fund ‘A way to achieve Europe’. The work was also  
558 supported by grants CSIC-COV19-014 from Consejo Superior de Investigaciones  
559 Científicas (CSIC), project 525/C/2021 from Fundació La Marató de TV3, PID2020-  
560 113888RB-I00 from Ministerio de Ciencia e Innovación, BFU2017-91384-EXP from  
561 Ministerio de Ciencia, Innovación y Universidades (MCIU), PI18/00210 and  
562 PI21/00139 from Instituto de Salud Carlos III and S2018/BAA-4370 (PLATESA2 from  
563 Comunidad de Madrid/FEDER). C.P., M.C. and P.M. are supported by the Miguel  
564 Servet programme of the Instituto de Salud Carlos III (CPII19/00001, CPII17/00006  
565 and CP16/00116, respectively) cofinanced by the European Regional Development  
566 Fund (ERDF). CIBERehd (Centro de Investigación en Red de Enfermedades Hepáticas  
567 y Digestivas) is funded by Instituto de Salud Carlos III. Institutional grants from the

568 Fundación Ramón Areces and Banco Santander to the CBMSO are also acknowledged.  
569 The team at CBMSO belongs to the Global Virus Network (GVN). B.M.-G. is  
570 supported by predoctoral contract PFIS FI19/00119 from Instituto de Salud Carlos III  
571 (Ministerio de Sanidad y Consumo) cofinanced by Fondo Social Europeo (FSE). R.L.-  
572 V. is supported by predoctoral contract PEJD-2019-PRE/BMD-16414 from Comunidad  
573 de Madrid. C.G.-C. is supported by predoctoral contract PRE2018-083422 from MCIU.  
574 BS was supported by a predoctoral research fellowship (Doctorados Industriales, DI-17-  
575 09134) from Spanish MINECO.

## 576 REFERENCES

- 577 1. Huang C, Wang Y, Li X, Ren L, Zhao J, Hu Y, Zhang L, Fan G, Xu J, Gu X,  
578 Cheng Z, Yu T, Xia J, Wei Y, Wu W, Xie X, Yin W, Li H, Liu M, Xiao Y, Gao  
579 H, Guo L, Xie J, Wang G, Jiang R, Gao Z, Jin Q, Wang J, Cao B. 2020. Clinical  
580 features of patients infected with 2019 novel coronavirus in Wuhan, China.  
581 *Lancet* 395:497-506.
- 582 2. Dos Santos WG. 2021. Impact of virus genetic variability and host immunity for  
583 the success of COVID-19 vaccines. *Biomed Pharmacother* 136:111272.
- 584 3. Tillett RL, Sevinsky JR, Hartley PD, Kerwin H, Crawford N, Gorzalski A,  
585 Laverdure C, Verma SC, Rossetto CC, Jackson D, Farrell MJ, Van Hooser S,  
586 Pandori M. 2021. Genomic evidence for reinfection with SARS-CoV-2: a case  
587 study. *Lancet Infect Dis* 21:52-58.
- 588 4. To KK, Hung IF, Ip JD, Chu AW, Chan WM, Tam AR, Fong CH, Yuan S, Tsoi  
589 HW, Ng AC, Lee LL, Wan P, Tso E, To WK, Tsang D, Chan KH, Huang JD,  
590 Kok KH, Cheng VC, Yuen KY. 2021. COVID-19 re-infection by a  
591 phylogenetically distinct SARS-coronavirus-2 strain confirmed by whole  
592 genome sequencing. *Clin Infect Dis* 73(9):e2946-e2951.
- 593 5. Lee M. 2021. Lack of Severe Acute Respiratory Syndrome Coronavirus 2  
594 Neutralization by Antibodies to Seasonal Coronaviruses: Making Sense of the  
595 Coronavirus Disease 2019 Pandemic. *Clin Infect Dis* 73:e1212-e1213.
- 596 6. Baum A, Fulton BO, Wloga E, Copin R, Pascal KE, Russo V, Giordano S,  
597 Lanza K, Negron N, Ni M, Wei Y, Atwal GS, Murphy AJ, Stahl N,  
598 Yancopoulos GD, Kyrtatsous CA. 2020. Antibody cocktail to SARS-CoV-2  
599 spike protein prevents rapid mutational escape seen with individual antibodies.  
600 *Science* 369:1014-1018.
- 601 7. Hacisuleyman E, Hale C, Saito Y, Blachere NE, Bergh M, Conlon EG,  
602 Schaefer-Babajew DJ, DaSilva J, Muecksch F, Gaebler C, Lifton R,  
603 Nussenzweig MC, Hatzioannou T, Bieniasz PD, Darnell RB. 2021. Vaccine  
604 Breakthrough Infections with SARS-CoV-2 Variants. *N Engl J Med* 384:2212-  
605 2218.
- 606 8. McCallum M, Bassi J, De Marco A, Chen A, Walls AC, Di Iulio J, Tortorici  
607 MA, Navarro MJ, Silacci-Fregni C, Saliba C, Sprouse KR, Agostini M, Pinto D,  
608 Culap K, Bianchi S, Jaconi S, Cameroni E, Bowen JE, Tilles SW, Pizzuto MS,  
609 Guastalla SB, Bona G, Pellanda AF, Garzoni C, Van Voorhis WC, Rosen LE,  
610 Snell G, Telenti A, Virgin HW, Piccoli L, Corti D, Velesler D. 2021. SARS-  
611 CoV-2 immune evasion by the B.1.427/B.1.429 variant of concern. *Science*  
612 373:648-654.

- 613 9. Nonaka CKV, Franco MM, Graf T, de Lorenzo Barcia CA, de Avila Mendonca  
614 RN, de Sousa KAF, Neiva LMC, Fosenca V, Mendes AVA, de Aguiar RS,  
615 Giovanetti M, de Freitas Souza BS. 2021. Genomic Evidence of SARS-CoV-2  
616 Reinfection Involving E484K Spike Mutation, Brazil. *Emerg Infect Dis*  
617 27:1522-1524.
- 618 10. Weisblum Y, Schmidt F, Zhang F, DaSilva J, Poston D, Lorenzi JC, Muecksch  
619 F, Rutkowska M, Hoffmann HH, Michailidis E, Gaebler C, Agudelo M, Cho A,  
620 Wang Z, Gazumyan A, Cipolla M, Luchsinger L, Hillyer CD, Caskey M,  
621 Robbiani DF, Rice CM, Nussenzweig MC, Hatzioannou T, Bieniasz PD. 2020.  
622 Escape from neutralizing antibodies by SARS-CoV-2 spike protein variants.  
623 *Elife* 9:e61312.
- 624 11. Domingo E, Perales C. 2019. Viral quasispecies. *PLoS Genet* 15:e1008271.
- 625 12. Domingo E, García-Crespo C, Perales C. 2021. Historical perspective on the  
626 discovery of the quasispecies concept. *Annu Rev Virol* 8:51-72.
- 627 13. Gregori J, Perales C, Rodriguez-Frias F, Esteban JI, Quer J, Domingo E. 2016.  
628 Viral quasispecies complexity measures. *Virology* 493:227-237.
- 629 14. Fuhrmann L, Jablonski KP, Beerenwinkel N. 2021. Quantitative measures of  
630 within-host viral genetic diversity. *Curr Opin Virol* 49:157-163.
- 631 15. Marcus PI, Rodriguez LL, Sekellick MJ. 1998. Interferon induction as a  
632 quasispecies marker of vesicular stomatitis virus populations. *J Virol* 72:542-  
633 549.
- 634 16. Farci P. 2001. Hepatitis C virus. The importance of viral heterogeneity. *Clin*  
635 *Liver Dis* 5:895-916.
- 636 17. Baranowski E, Ruiz-Jarabo CM, Pariente N, Verdaguer N, Domingo E. 2003.  
637 Evolution of cell recognition by viruses: a source of biological novelty with  
638 medical implications. *Adv Virus Res* 62:19-111.
- 639 18. Farci P. 2011. New insights into the HCV quasispecies and  
640 compartmentalization. *Semin Liver Dis* 31:356-74.
- 641 19. Domingo E, Sheldon J, Perales C. 2012. Viral quasispecies evolution. *Microbiol*  
642 *Mol Biol Rev* 76:159-216.
- 643 20. Honce R, Schultz-Cherry S. 2020. They are what you eat: Shaping of viral  
644 populations through nutrition and consequences for virulence. *PLoS Pathog*  
645 16:e1008711.
- 646 21. Young DF, Wignall-Fleming EB, Busse DC, Pickin MJ, Hankinson J, Randall  
647 EM, Tavendale A, Davison AJ, Lamont D, Tregoning JS, Goodbourn S, Randall  
648 RE. 2019. The switch between acute and persistent paramyxovirus infection  
649 caused by single amino acid substitutions in the RNA polymerase P subunit.  
650 *PLoS Pathog* 15:e1007561.
- 651 22. Rima BK, Gatherer D, Young DF, Norsted H, Randall RE, Davison AJ. 2014.  
652 Stability of the parainfluenza virus 5 genome revealed by deep sequencing of  
653 strains isolated from different hosts and following passage in cell culture. *J Virol*  
654 88:3826-36.
- 655 23. Karamitros T, Papadopoulou G, Bousali M, Mexias A, Tsiodras S, Mentis A.  
656 2020. SARS-CoV-2 exhibits intra-host genomic plasticity and low-frequency  
657 polymorphic quasispecies. *J Clin Virol* 131:104585.
- 658 24. Jary A, Leducq V, Malet I, Marot S, Klement-Frutos E, Teyssou E, Soulie C,  
659 Abdi B, Wirden M, Pourcher V, Caumes E, Calvez V, Burrel S, Marcelin AG,  
660 Boutolleau D. 2020. Evolution of viral quasispecies during SARS-CoV-2  
661 infection. *Clin Microbiol Infect* 26(11):1560.e1-1560.e4.

- 662 25. Rueca M, Bartolini B, Gruber CEM, Piralla A, Baldanti F, Giombini E, Messina  
663 F, Marchioni L, Ippolito G, Di Caro A, Capobianchi MR. 2020.  
664 Compartmentalized Replication of SARS-Cov-2 in Upper vs. Lower Respiratory  
665 Tract Assessed by Whole Genome Quasispecies Analysis. *Microorganisms*  
666 8:E1302.
- 667 26. Capobianchi MR, Rueca M, Messina F, Giombini E, Carletti F, Colavita F,  
668 Castilletti C, Lalle E, Bordi L, Vairo F, Nicastrì E, Ippolito G, Gruber CEM,  
669 Bartolini B. 2020. Molecular characterization of SARS-CoV-2 from the first  
670 case of COVID-19 in Italy. *Clin Microbiol Infect* 26:954-956.
- 671 27. Sun F, Wang X, Tan S, Dan Y, Lu Y, Zhang J, Xu J, Tan Z, Xiang X, Zhou Y,  
672 He W, Wan X, Zhang W, Chen Y, Tan W, Deng G. 2021. SARS-CoV-2  
673 Quasispecies Provides an Advantage Mutation Pool for the Epidemic Variants.  
674 *Microbiol Spectr* 9:e0026121.
- 675 28. Andres C, Garcia-Cehic D, Gregori J, Pinana M, Rodriguez-Frias F, Guerrero-  
676 Murillo M, Esperalba J, Rando A, Goterris L, Codina MG, Quer S, Martin MC,  
677 Campins M, Ferrer R, Almirante B, Esteban JI, Pumarola T, Anton A, Quer J.  
678 2020. Naturally occurring SARS-CoV-2 gene deletions close to the spike S1/S2  
679 cleavage site in the viral quasispecies of COVID19 patients. *Emerg Microbes*  
680 *Infect* 9:1900-1911.
- 681 29. Wong YC, Lau SY, Wang To KK, Mok BWY, Li X, Wang P, Deng S, Woo KF,  
682 Du Z, Li C, Zhou J, Chan JFW, Yuen KY, Chen H, Chen Z. 2021. Natural  
683 Transmission of Bat-like Severe Acute Respiratory Syndrome Coronavirus 2  
684 Without Proline-Arginine-Arginine-Alanine Variants in Coronavirus Disease  
685 2019 Patients. *Clin Infect Dis* 73:e437-e444.
- 686 30. Xu D, Zhang Z, Wang FS. 2004. SARS-associated coronavirus quasispecies in  
687 individual patients. *N Engl J Med* 350:1366-7.
- 688 31. Tang JW, Cheung JL, Chu IM, Sung JJ, Peiris M, Chan PK. 2006. The large  
689 386-nt deletion in SARS-associated coronavirus: evidence for quasispecies? *J*  
690 *Infect Dis* 194:808-13.
- 691 32. Liu J, Lim SL, Ruan Y, Ling AE, Ng LF, Drosten C, Liu ET, Stanton LW,  
692 Hibberd ML. 2005. SARS transmission pattern in Singapore reassessed by viral  
693 sequence variation analysis. *PLoS Med* 2:e43.
- 694 33. Park D, Huh HJ, Kim YJ, Son DS, Jeon HJ, Im EH, Kim JW, Lee NY, Kang ES,  
695 Kang CI, Chung DR, Ahn JH, Peck KR, Choi SS, Kim YJ, Ki CS, Park WY.  
696 2016. Analysis of inpatient heterogeneity uncovers the microevolution of  
697 Middle East respiratory syndrome coronavirus. *Cold Spring Harb Mol Case Stud*  
698 2:a001214.
- 699 34. Borucki MK, Lao V, Hwang M, Gardner S, Adney D, Munster V, Bowen R,  
700 Allen JE. 2016. Middle East Respiratory Syndrome Coronavirus Intra-Host  
701 Populations Are Characterized by Numerous High Frequency Variants. *PLoS*  
702 *One* 11:e0146251.
- 703 35. Gregori J, Cortese MF, Pinana M, Campos C, Garcia-Cehic D, Andres C, Abril  
704 JF, Codina MG, Rando A, Esperalba J, Sulleiro E, Joseph J, Saubi N, Colomer-  
705 Castell S, Martin MC, Castillo C, Esteban JI, Pumarola T, Rodriguez-Frias F,  
706 Anton A, Quer J. 2021. Host-dependent editing of SARS-CoV-2 in COVID-19  
707 patients. *Emerg Microbes Infect* 10:1777-1789.
- 708 36. Al Khatib HA, Benslimane FM, Elbashir IE, Coyle PV, Al Maslamani MA, Al-  
709 Khal A, Al Thani AA, Yassine HM. 2020. Within-Host Diversity of SARS-  
710 CoV-2 in COVID-19 Patients With Variable Disease Severities. *Front Cell*  
711 *Infect Microbiol* 10:575613.

- 712 37. Soria ME, Corton M, Martinez-Gonzalez B, Lobo-Vega R, Vazquez-Sirvent L,  
713 Lopez-Rodriguez R, Almoguera B, Mahillo I, Minguez P, Herrero A, Taracido  
714 JC, Macias-Valcayo A, Esteban J, Fernandez-Roblas R, Gadea I, Ruiz-Hornillos  
715 J, Ayuso C, Perales C. 2021. High SARS-CoV-2 viral load is associated with a  
716 worse clinical outcome of COVID-19 disease. *Access Microbiol* 3:000259.
- 717 38. Soria ME, Gregori J, Chen Q, Garcia-Cehic D, Llorens M, de Avila AI, Beach  
718 NM, Domingo E, Rodriguez-Frias F, Buti M, Esteban R, Esteban JI, Quer J,  
719 Perales C. 2018. Pipeline for specific subtype amplification and drug resistance  
720 detection in hepatitis C virus. *BMC Infect Dis* 18:446.
- 721 39. Soria ME, Garcia-Crespo C, Martinez-Gonzalez B, Vazquez-Sirvent L, Lobo-  
722 Vega R, de Avila AI, Gallego I, Chen Q, Garcia-Cehic D, Llorens-Revull M,  
723 Briones C, Gomez J, Ferrer-Orta C, Verdaguer N, Gregori J, Rodriguez-Frias F,  
724 Buti M, Esteban JI, Domingo E, Quer J, Perales C. 2020. Amino Acid  
725 Substitutions Associated with Treatment Failure for Hepatitis C Virus Infection.  
726 *J Clin Microbiol* 58:e01985-20.
- 727 40. Chen Q, Perales C, Soria ME, Garcia-Cehic D, Gregori J, Rodriguez-Frias F,  
728 Buti M, Crespo J, Calleja JL, Tabernero D, Vila M, Lazaro F, Rando-Segura A,  
729 Nieto-Aponte L, Llorens-Revull M, Cortese MF, Fernandez-Alonso I, Castellote  
730 J, Niubo J, Imaz A, Xiol X, Castells L, Riveiro-Barciela M, Llaneras J, Navarro  
731 J, Vargas-Blasco V, Augustin S, Conde I, Rubin A, Prieto M, Torras X, Margall  
732 N, Forns X, Marino Z, Lens S, Bonacci M, Perez-Del-Pulgar S, Londono MC,  
733 Garcia-Buey ML, Sanz-Cameno P, Morillas R, Martro E, Saludes V, Masnou-  
734 Ridaura H, Salmeron J, Quiles R, Carrion JA, Forne M, Rosinach M, Fernandez  
735 I, et al. 2020. Deep-sequencing reveals broad subtype-specific HCV resistance  
736 mutations associated with treatment failure. *Antiviral Res* 174:104694.
- 737 41. Perales C, Chen Q, Soria ME, Gregori J, Garcia-Cehic D, Nieto-Aponte L,  
738 Castells L, Imaz A, Llorens-Revull M, Domingo E, Buti M, Esteban JI,  
739 Rodriguez-Frias F, Quer J. 2018. Baseline hepatitis C virus resistance-associated  
740 substitutions present at frequencies lower than 15% may be clinically  
741 significant. *Infect Drug Resist* 11:2207-2210.
- 742 42. Hathaway NJ, Parobek CM, Juliano JJ, Bailey JA. 2018. SeekDeep: single-base  
743 resolution de novo clustering for amplicon deep sequencing. *Nucleic Acids Res*  
744 46:e21.
- 745 43. Guerrero-Murillo M, Gregori i Font J. 2018. QSutils: Quasispecies Diversity. R  
746 package version 1.0.0.
- 747 44. Feng DF, Doolittle RF. 1996. Progressive alignment of amino acid sequences  
748 and construction of phylogenetic trees from them. *Methods in Enzymol*  
749 266:368-82.
- 750 45. Hecht M, Bromberg Y, Rost B. 2015. Better prediction of functional effects for  
751 sequence variants. *BMC Genomics* 16 Suppl 8:S1.
- 752 46. Alouane T, Laamarti M, Essabbar A, Hakmi M, Bouricha EM, Chemaou-Elfihri  
753 MW, Kartti S, Boumajdi N, Bendani H, Laamarti R, Ghrifi F, Allam L, Aanniz  
754 T, Ouadghiri M, El Hafidi N, El Jaoudi R, Benrahma H, Attar JE, Mentag R,  
755 Sbabou L, Nejari C, Amzazi S, Belyamani L, Ibrahimi A. 2020. Genomic  
756 Diversity and Hotspot Mutations in 30,983 SARS-CoV-2 Genomes: Moving  
757 Toward a Universal Vaccine for the "Confined Virus"? *Pathogens* 9(10):829.
- 758 47. Badua C, Baldo KAT, Medina PMB. 2021. Genomic and proteomic mutation  
759 landscapes of SARS-CoV-2. *J Med Virol* 93:1702-1721.
- 760 48. Laamarti M, Alouane T, Kartti S, Chemaou-Elfihri MW, Hakmi M, Essabbar A,  
761 Laamarti M, Hlali H, Bendani H, Boumajdi N, Benhrif O, Allam L, El Hafidi N,

- 762 El Jaoudi R, Allali I, Marchoudi N, Fekak J, Benrahma H, Nejjari C, Amzazi S,  
763 Belyamani L, Ibrahimi A. 2020. Large scale genomic analysis of 3067 SARS-  
764 CoV-2 genomes reveals a clonal geo-distribution and a rich genetic variations of  
765 hotspots mutations. *PLoS One* 15:e0240345.
- 766 49. Shu Y, McCauley J. 2017. GISAID: Global initiative on sharing all influenza  
767 data - from vision to reality. *Euro Surveill* 22:30494.
- 768 50. Mourier T, Sadykov M, Carr MJ, Gonzalez G, Hall WW, Pain A. 2021. Host-  
769 directed editing of the SARS-CoV-2 genome. *Biochem Biophys Res Commun*  
770 538:35-39.
- 771 51. Smith EC, Case JB, Blanc H, Isakov O, Shomron N, Vignuzzi M, Denison MR.  
772 2015. Mutations in coronavirus nonstructural protein 10 decrease virus  
773 replication fidelity. *J Virol* 89:6418-26.
- 774 52. Stapleford KA, Rozen-Gagnon K, Das PK, Saul S, Poirier EZ, Blanc H,  
775 Vidalain PO, Merits A, Vignuzzi M. 2015. Viral Polymerase-Helicase  
776 Complexes Regulate Replication Fidelity To Overcome Intracellular Nucleotide  
777 Depletion. *J Virol* 89:11233-44.
- 778 53. Agudo R, de la Higuera I, Arias A, Grande-Perez A, Domingo E. 2016.  
779 Involvement of a joker mutation in a polymerase-independent lethal mutagenesis  
780 escape mechanism. *Virology* 494:257-266.
- 781 54. Collins ND, Beck AS, Widen SG, Wood TG, Higgs S, Barrett ADT. 2018.  
782 Structural and Nonstructural Genes Contribute to the Genetic Diversity of RNA  
783 Viruses. *mBio* 9:e01871-18.
- 784 55. V'Kovski P, Kratzel A, Steiner S, Stalder H, Thiel V. 2021. Coronavirus biology  
785 and replication: implications for SARS-CoV-2. *Nat Rev Microbiol* 19:155-170.
- 786 56. Cheemarla NR, Watkins TA, Mihaylova VT, Wang B, Zhao D, Wang G, Landry  
787 ML, Foxman EF. 2021. Dynamic innate immune response determines  
788 susceptibility to SARS-CoV-2 infection and early replication kinetics. *J Exp*  
789 *Med* 218:e20210583.
- 790 57. Harrison AG, Lin T, Wang P. 2020. Mechanisms of SARS-CoV-2 Transmission  
791 and Pathogenesis. *Trends Immunol* 41:1100-1115.
- 792 58. Farci P, Shimoda A, Coiana A, Diaz G, Peddis G, Melpolder JC, Strazzer A,  
793 Chien DY, Munoz SJ, Balestrieri A, Purcell RH, Alter HJ. 2000. The outcome  
794 of acute hepatitis C predicted by the evolution of the viral quasispecies. *Science*  
795 288:339-44.
- 796 59. Farci P, Strazzer R, Alter HJ, Farci S, Degioannis D, Coiana A, Peddis G, Usai  
797 F, Serra G, Chessa L, Diaz G, Balestrieri A, Purcell RH. 2002. Early changes in  
798 hepatitis C viral quasispecies during interferon therapy predict the therapeutic  
799 outcome. *Proc Natl Acad Sci U S A* 99:3081-6.
- 800 60. Li Q, Wu J, Nie J, Zhang L, Hao H, Liu S, Zhao C, Zhang Q, Liu H, Nie L, Qin  
801 H, Wang M, Lu Q, Li X, Sun Q, Liu J, Zhang L, Li X, Huang W, Wang Y.  
802 2020. The Impact of Mutations in SARS-CoV-2 Spike on Viral Infectivity and  
803 Antigenicity. *Cell* 182:1284-1294 e9.
- 804 61. Grabowski F, Preibisch G, Gizinski S, Kochanczyk M, Lipniacki T. 2021.  
805 SARS-CoV-2 Variant of Concern 202012/01 Has about Twofold Replicative  
806 Advantage and Acquires Concerning Mutations. *Viruses* 13:392.
- 807 62. Alenquer M, Ferreira F, Lousa D, Valerio M, Medina-Lopes M, Bergman ML,  
808 Goncalves J, Demengeot J, Leite RB, Lilue J, Ning Z, Penha-Goncalves C,  
809 Soares H, Soares CM, Amorim MJ. 2021. Signatures in SARS-CoV-2 spike  
810 protein conferring escape to neutralizing antibodies. *PLoS Pathog* 17:e1009772.

- 811 63. Wang R, Chen J, Gao K, Wei GW. 2021. Vaccine-escape and fast-growing  
812 mutations in the United Kingdom, the United States, Singapore, Spain, India,  
813 and other COVID-19-devastated countries. *Genomics* 113:2158-2170.
- 814 64. Sola I, Almazan F, Zuniga S, Enjuanes L. 2015. Continuous and Discontinuous  
815 RNA Synthesis in Coronaviruses. *Annu Rev Virol* 2:265-88.
- 816 65. Vignuzzi M, Lopez CB. 2019. Defective viral genomes are key drivers of the  
817 virus-host interaction. *Nat Microbiol* 4:1075-1087.
- 818 66. Posthuma CC, Te Velthuis AJW, Snijder EJ. 2017. Nidovirus RNA  
819 polymerases: Complex enzymes handling exceptional RNA genomes. *Virus Res*  
820 234:58-73.
- 821 67. Hillen HS, Kokic G, Farnung L, Dienemann C, Tegunov D, Cramer P. 2020.  
822 Structure of replicating SARS-CoV-2 polymerase. *Nature* 584:154-156.
- 823 68. Armero A, Berthet N, Avarre JC. 2021. Intra-Host Diversity of SARS-Cov-2  
824 Should Not Be Neglected: Case of the State of Victoria, Australia. *Viruses*  
825 13:133.
- 826 69. Domingo E, Schuster P. 2016. Quasispecies: from theory to experimental  
827 systems. *Current Topics in Microbiology and Immunology*. Vol. 392. Springer.
- 828 70. Garcia-Crespo C, Soria ME, Gallego I, Avila AI, Martinez-Gonzalez B,  
829 Vazquez-Sirvent L, Gomez J, Briones C, Gregori J, Quer J, Perales C, Domingo  
830 E. 2020. Dissimilar Conservation Pattern in Hepatitis C Virus Mutant Spectra,  
831 Consensus Sequences, and Data Banks. *J Clin Med* 9:3450.
- 832 71. Domingo E, Perales C. 2012. From quasispecies theory to viral quasispecies:  
833 how complexity has permeated virology. *Math Model Nat Phenom* 7:32-49.
- 834 72. Sarkar R, Mitra S, Chandra P, Saha P, Banerjee A, Dutta S, Chawla-Sarkar M.  
835 2021. Comprehensive analysis of genomic diversity of SARS-CoV-2 in different  
836 geographic regions of India: an endeavour to classify Indian SARS-CoV-2  
837 strains on the basis of co-existing mutations. *Arch Virol* 166:801-812.
- 838 73. Simmonds P. 2020. Rampant C-->U Hypermutation in the Genomes of SARS-  
839 CoV-2 and Other Coronaviruses: Causes and Consequences for Their Short- and  
840 Long-Term Evolutionary Trajectories. *mSphere* 5:e00408-20.
- 841 74. Danchin A, Marliere P. 2020. Cytosine drives evolution of SARS-CoV-2.  
842 *Environ Microbiol* 22:1977-1985.
- 843 75. Li X, Zai J, Zhao Q, Nie Q, Li Y, Foley BT, Chaillon A. 2020. Evolutionary  
844 history, potential intermediate animal host, and cross-species analyses of SARS-  
845 CoV-2. *J Med Virol* 92:602-611.
- 846 76. Nie Q, Li X, Chen W, Liu D, Chen Y, Li H, Li D, Tian M, Tan W, Zai J. 2020.  
847 Phylogenetic and phylodynamic analyses of SARS-CoV-2. *Virus Res*  
848 287:198098.
- 849 77. Bai Y, Jiang D, Lon JR, Chen X, Hu M, Lin S, Chen Z, Wang X, Meng Y, Du  
850 H. 2020. Comprehensive evolution and molecular characteristics of a large  
851 number of SARS-CoV-2 genomes reveal its epidemic trends. *Int J Infect Dis*  
852 100:164-173.
- 853 78. Lai A, Bergna A, Acciarri C, Galli M, Zehender G. 2020. Early phylogenetic  
854 estimate of the effective reproduction number of SARS-CoV-2. *J Med Virol*  
855 92:675-679.
- 856 79. Nabil B, Sabrina B, Abdelhakim B. 2021. Transmission route and introduction  
857 of pandemic SARS-CoV-2 between China, Italy, and Spain. *J Med Virol*  
858 93:564-568.

- 859 80. Pereson MJ, Mojsiejczuk L, Martinez AP, Flichman DM, Garcia GH, Di Lello  
860 FA. 2021. Phylogenetic analysis of SARS-CoV-2 in the first few months since  
861 its emergence. *J Med Virol* 93:1722-1731.
- 862 81. Castells M, Lopez-Tort F, Colina R, Cristina J. 2020. Evidence of increasing  
863 diversification of emerging Severe Acute Respiratory Syndrome Coronavirus 2  
864 strains. *J Med Virol* 92:2165-2172.
- 865 82. Diez-Fuertes F, Iglesias-Caballero M, Garcia-Perez J, Monzon S, Jimenez P,  
866 Varona S, Cuesta I, Zaballos A, Jimenez M, Checa L, Pozo F, Perez-Olmeda M,  
867 Thomson MM, Alcami J, Casas I. 2021. A Founder Effect Led Early SARS-  
868 CoV-2 Transmission in Spain. *J Virol* 95:e01583-20.
- 869 83. Liu Q, Zhao S, Shi CM, Song S, Zhu S, Su Y, Zhao W, Li M, Bao Y, Xue Y,  
870 Chen H. 2020. Population Genetics of SARS-CoV-2: Disentangling Effects of  
871 Sampling Bias and Infection Clusters. *Genomics Proteomics Bioinformatics*  
872 18:640-647.
- 873 84. Domingo E. 2020. *Virus as Populations*. Academic Press, Elsevier, Amsterdam.  
874 Second Edition.
- 875 85. Domingo E, Garcia-Crespo C, Lobo-Vega R, Perales C. 2021. Mutation Rates,  
876 Mutation Frequencies, and Proofreading-Repair Activities in RNA Virus  
877 Genetics. *Viruses* 13:1882.
- 878 86. Wu F, Zhao S, Yu B, Chen YM, Wang W, Song ZG, Hu Y, Tao ZW, Tian JH,  
879 Pei YY, Yuan ML, Zhang YL, Dai FH, Liu Y, Wang QM, Zheng JJ, Xu L,  
880 Holmes EC, Zhang YZ. 2020. A new coronavirus associated with human  
881 respiratory disease in China. *Nature* 579:265-269.
- 882 87. de Avila AI, Gallego I, Soria ME, Gregori J, Quer J, Esteban JI, Rice CM,  
883 Domingo E, Perales C. 2016. Lethal Mutagenesis of Hepatitis C Virus Induced  
884 by Favipiravir. *PLoS ONE* 11:e0164691.
- 885 88. Gallego I, Soria ME, Gregori J, de Avila AI, Garcia-Crespo C, Moreno E, Gadea  
886 I, Esteban J, Fernandez-Roblas R, Esteban JI, Gomez J, Quer J, Domingo E,  
887 Perales C. 2019. Synergistic lethal mutagenesis of hepatitis C virus. *Antimicrob*  
888 *Agents Chemother* 63:e01653-19.
- 889 89. García-Crespo C, Gallego I, Soria ME, De Ávila AI, Martínez-González B,  
890 Vázquez-Sirvent L, Lobo-Vega R, Moreno E, Gómez J, Briones C, Gregori J,  
891 Quer J, Domingo E, Perales C. 2021. Population disequilibrium as promoter of  
892 adaptive explorations in hepatitis C virus. *Viruses* 13:616.
- 893 90. Gregori J, Salicru M, Domingo E, Sanchez A, Esteban JI, Rodriguez-Frias F,  
894 Quer J. 2014. Inference with viral quasispecies diversity indices: clonal and  
895 NGS approaches. *Bioinformatics* 30:1104-1111.

896

## 897 **FIGURE LEGENDS**

898 **FIG 1.** Representation of the SARS-CoV-2 genome, encoded proteins, and amplicons  
899 analyzed by UDS. The region corresponding to the ORF1ab of the virus is shown at the  
900 top. In the two boxes, the nsp12 (blue) and spike (orange) have been expanded, with the  
901 first and last nucleotide number given at the beginning and the end of the bars,  
902 respectively (genome numbering is according to the reference genome NCBI accession  
903 number: NC\_045512.2). Relevant protein domains are indicated, including motifs A to



904 G depicted as protruding grey boxes in nsp12 (polymerase), and the receptor binding  
905 motif (RBM) and the S1/S2 cleavage site in S. The amplicons [A1 to A4 for the nsp12  
906 (polymerase) and A5, A6 for S] are shown flanked by horizontal arrows that mark the  
907 position of the oligonucleotide primers used for amplification (oligonucleotide  
908 sequences are given in Table S5 in <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>).  
909 Flanking black boxes indicate the amino acids (aa) of nsp12 (polymerase) and S  
910 covered by the amplicons.

911 **FIG 2.** Heat map of point mutation and deletion frequencies in mutant spectra of SARS-  
912 CoV-2 from individual patients. Data are presented in two blocks, one for the nsp12  
913 (polymerase)-coding region (genomic residues 14,534-16,054), and another for the S-  
914 coding region (genomic residues 22,872-23,645). Only positions with a mutation or  
915 those affected by a deletion are represented. Each row corresponds to a patient, and  
916 patients have been divided in those with mild, moderate and exitus disease outcomes  
917 (color coded, and with the patient identification code written at the left of each row).  
918 The patients' clinical status and demographic data are described in Table S1 in  
919 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>. Mutations and deletions have been  
920 identified relative to NCBI reference sequence NC\_045512.2. Each mutation and  
921 deletion (delta symbol  $\Delta$ ) with a frequency above the cut-off level (0.5%) is indicated,  
922 and its frequency within the mutant spectrum retrieved from each patient has been  
923 visualized with a color code displayed in the heading boxes (top left of the two blocs).  
924 Procedures are detailed in Materials and Methods.

925 **FIG 3.** Point mutations and deletions in the mutant spectra of SARS-CoV-2 isolates,  
926 distributed according to COVID-19 severity. The point mutations and deletions are  
927 those depicted in Fig. 2. **(A)** Total number of different point mutations in the nsp12  
928 (polymerase)- (left panel) and the S- (right panel) coding region distributed according to  
929 disease severity (mild, moderate, exitus, as indicated in the abscissa) in the patients  
930 from whom the virus was isolated. Bars indicate the total absolute number of mutations  
931 (left ordinate axes) and empty dots give the percentages normalized to the length in  
932 nucleotides of the sequenced regions (right ordinate axes). **(B)** Total number of different  
933 deletions in the nsp12 (polymerase)- (left panel) and the S- (right panel) coding region  
934 distributed according to disease severity in the patients from whom the virus was  
935 isolated. For **(A)** and **(B)** the statistical significance of the differences was determined  
936 by the proportion test; ns; not significant, \*\*\* $p < 0.001$ .

937 **FIG 4.** Comparison of the diversity indices for all amplicons of either the nsp12  
938 (polymerase)- or S-coding region, distributed according to virus-associated disease  
939 severity. The types of indices (abundance or incidence) are indicated in the heading  
940 filled boxes. The specific index is indicated in ordinate (13) ( $H_s$ , Shannon entropy;  
941  $Mf_{max}$ , maximum mutation frequency;  $H_{GS}$ , Gini Simpson;  $\pi$ , nucleotide diversity; N.  
942 poly.sites, number of polymorphic sites; N. hpl., number of haplotypes). Each cross is  
943 the numerical value obtained for the virus of an individual patient; patients have been  
944 distributed according to disease severity as indicated in abscissa (color coded). Data  
945 were obtained using a cut-off value of 0.1%, as previously reported (13). Values for  
946 each amplicon and patient are compiled in Table S3 in  
947 <https://saco.csic.es/index.php/s/8GH5aJgritCjEx5>. The statistical significance of the  
948 differences has been determined by the Wilcoxon test. \*,  $p < 0.05$ ; \*\*,  $p < 0.01$ ;  
949 \*\*\* $p < 0.001$ ; absence of connecting lines means that the difference between two patient  
950 groups was not statistically significant.

951

952 **FIG 5.** Location of amino acid substitutions in the three-dimensional structure of nsp12  
953 (polymerase). The structure used as reference is that of the replication complex nsp12-  
954 nsp8-nsp7 (PDB code 6NUR with the RNA superimposed from 7CYQ). (A)  
955 Substitutions found at low frequency (0.5% to 2%) in the mutant spectra. The central  
956 structure is a cartoon representation of the nsp12, depicted in grey and green, the latter  
957 showing the regions covered by amplicons A1-A4 (indicated in Fig. 1). Contact proteins  
958 nsp8 (orange) and nsp7 (yellow) are also drawn. Substitutions are labeled, and amino  
959 acids are shown as sticks in different colors, according to associated disease category:  
960 exitus in red; mild disease in yellow. Insets highlight the interactions of some  
961 substitutions with neighboring residues within a 5-Å radius. Two insets are shown per  
962 position, indicating the original and mutated residues, squared in blue and red,  
963 respectively. (B) Same design as A but with substitutions found at frequency higher  
964 than 2%. The substitutions, their frequency in the mutant spectrum, acceptability,  
965 functional score, and possible structural or functional effects are listed in Table 2.

966

967 **FIG 6.** Location of amino acid substitutions in the three-dimensional structure of spike  
968 (S) protein. The central structure is a cartoon representation of S trimer (PDB code  
969 7A94) with the reference monomer colored in green and dark-green, the latter marking

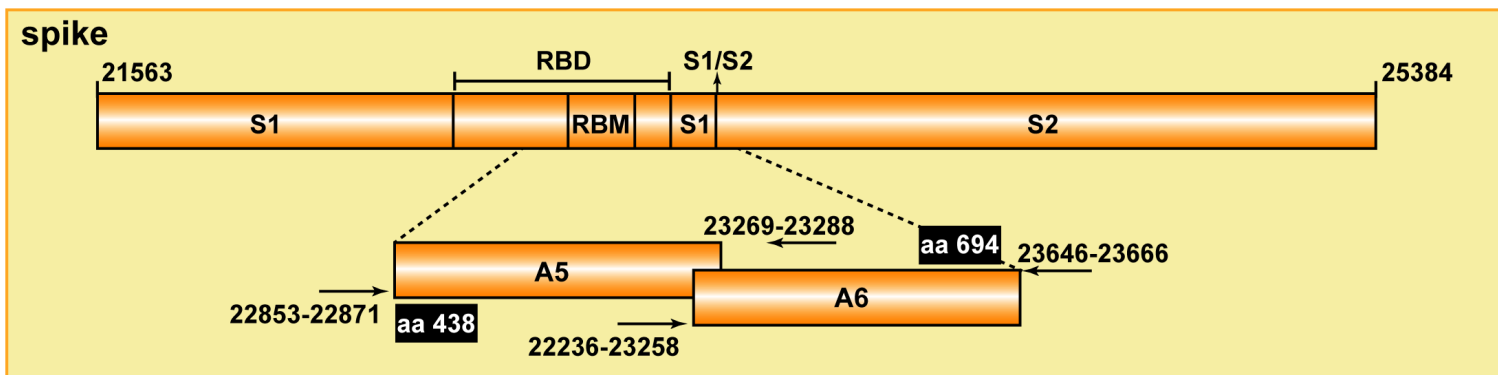
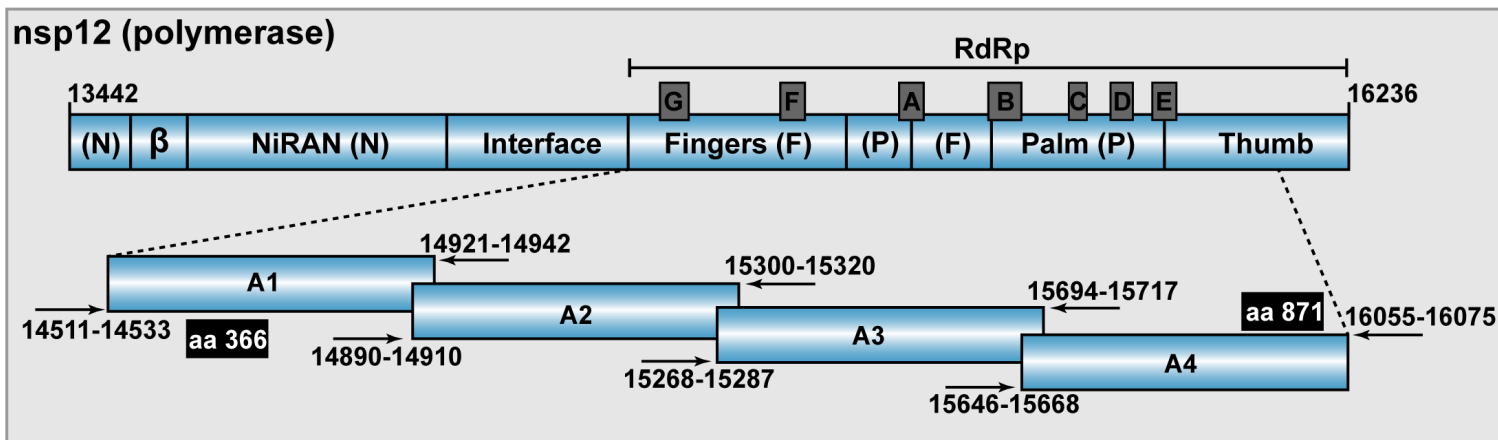
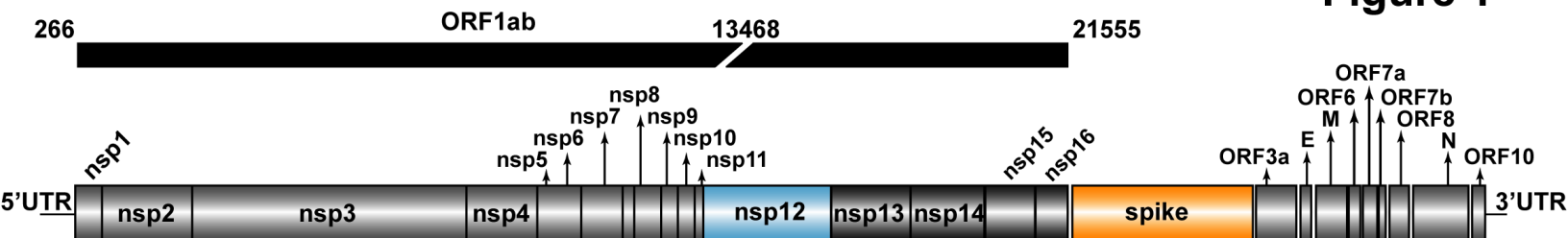
970 the regions covered by amplicons A5-A6 (indicated in Fig. 1). The remaining  
971 monomers of the S trimer are shown in grey. The reference monomer contains de RBD  
972 domain in the “erect” position. A superimposition of this domain in the “open”  
973 conformation is also shown in orange. Substitutions are labeled, and amino acids are  
974 shown as sticks in different colors, according to associated disease category: exitus in  
975 red; moderate in magenta and mild in yellow. Insets highlight the interactions of some  
976 substituted positions with neighboring residues within a 5-Å radius. Except for position  
977 522, two insets are shown per mutated position, indicating the original and mutated  
978 residues, squared in blue and red, respectively. For position 522, four insets are shown;  
979 the top two indicate the interactions of this residue in the open conformation of RBD,  
980 and the bottom two in the erect conformation. The substitutions, their frequency in the  
981 mutant spectrum, acceptability, functional score, and possible structural or functional  
982 effects are listed in Table 3.

983 **FIG 7.** Point mutation and deletion hot spots in SARS-CoV-2 mutant spectra. **(A)**  
984 Distribution of the total number of different variations (point mutations and deletions,  
985 given in ordinate) among the amplicons analyzed (indicated in abscissa). The statistical  
986 significance of the differences was determined by the proportion test. \*\*\*,  $p < 0.001$ ;  
987 absence of connecting lines among nsp12 amplicons means that differences were not  
988 statistically significant. **(B)** Location of point mutations and deletions within each  
989 amplicon (indicated in each box). Genome residue numbering is according to reference  
990 NCBI accession number NC\_045512.2. The numbers written in a yellow box refer to  
991 the number of patients whose virus carried the same mutation or deletion, and serve to  
992 identify hot spots. Point mutations and deletions were counted relative to the consensus  
993 sequence of the corresponding population.

994

995

Figure 1



**Figure 2**

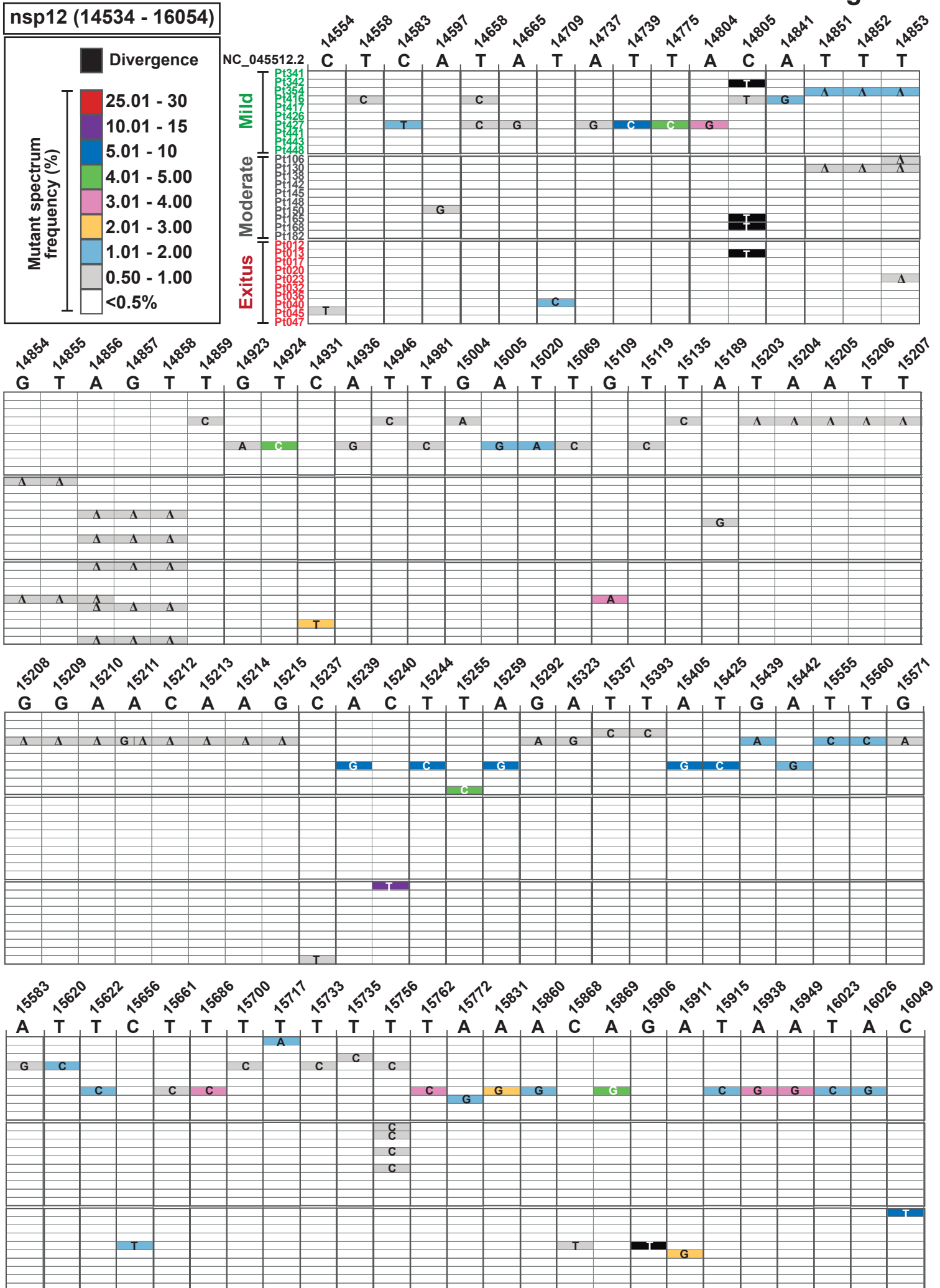




Figure 3

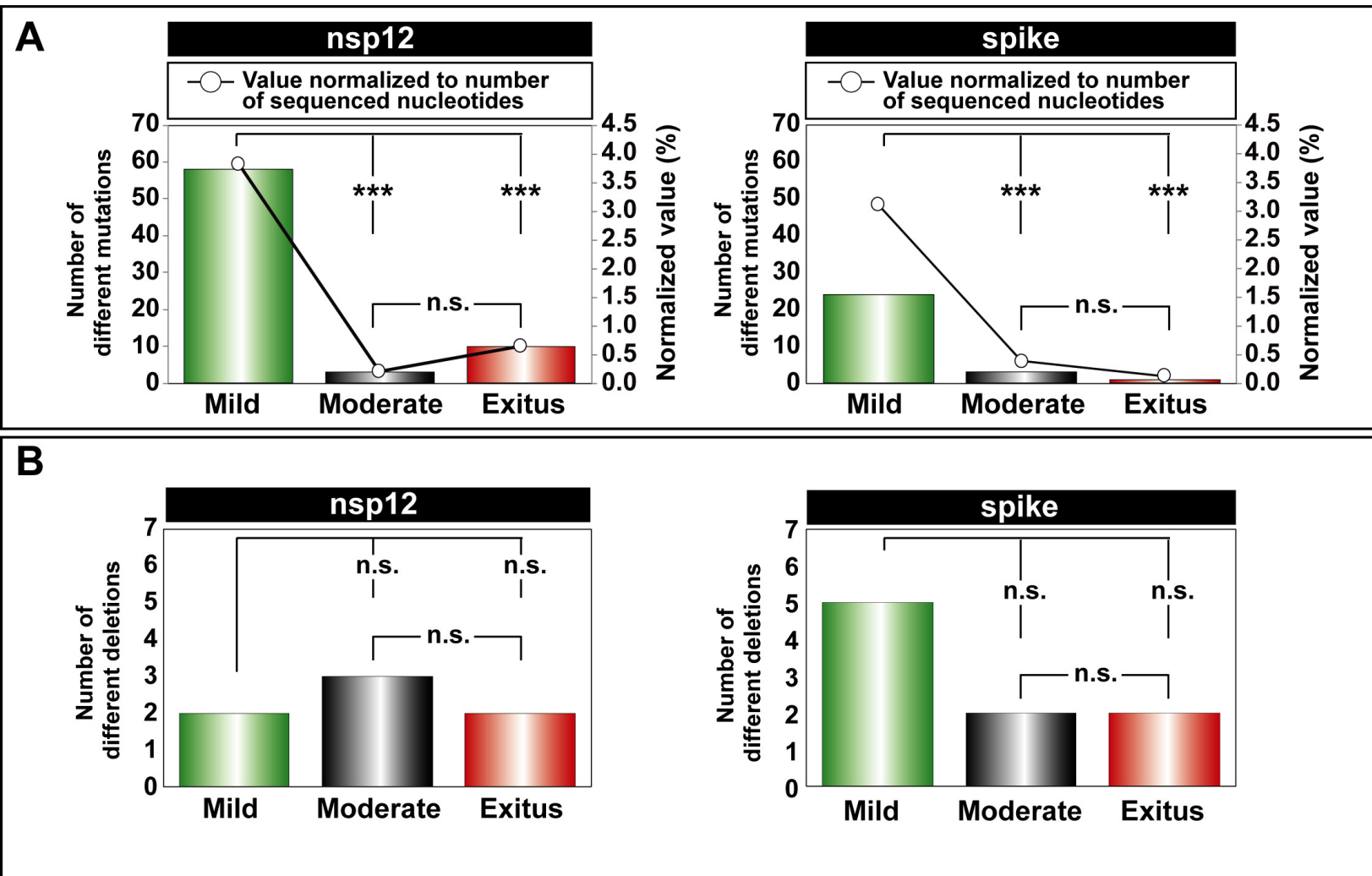
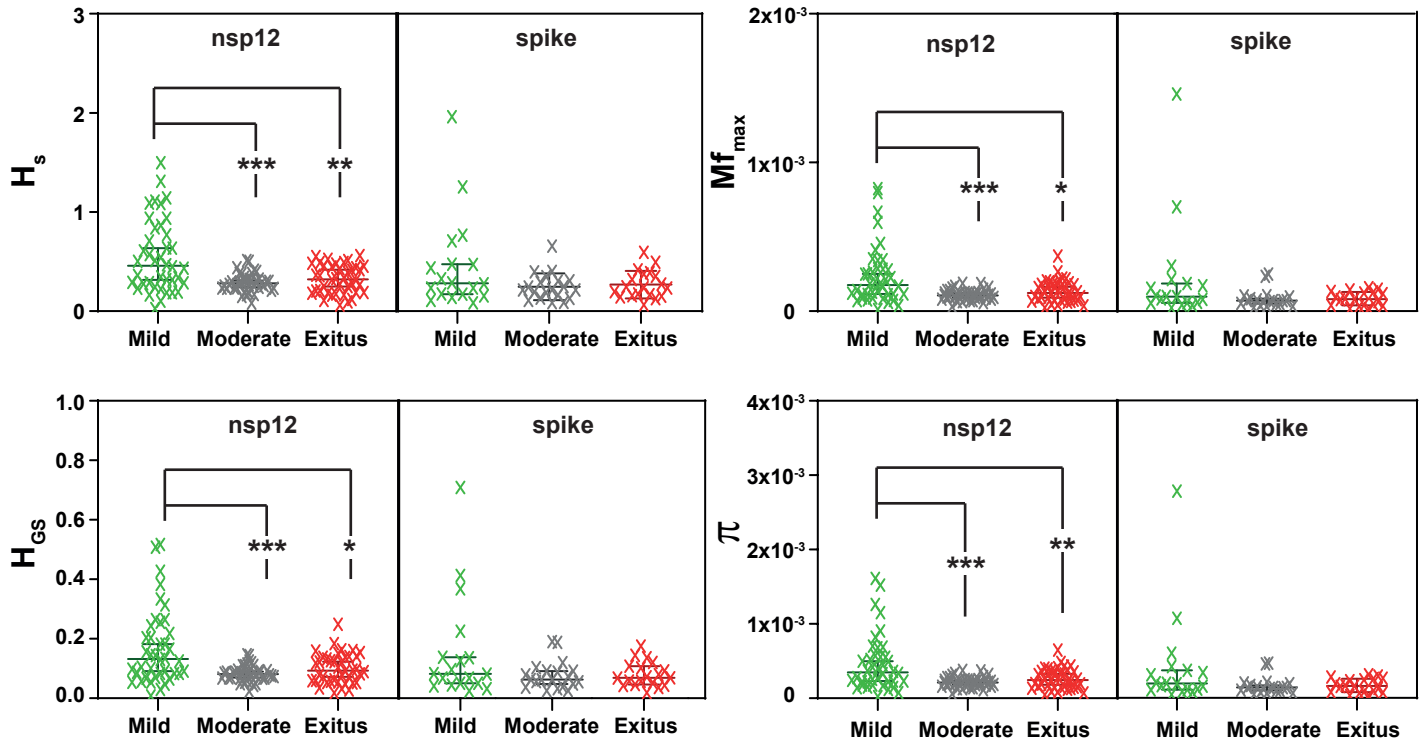


Figure 4

Abundance



Incidence

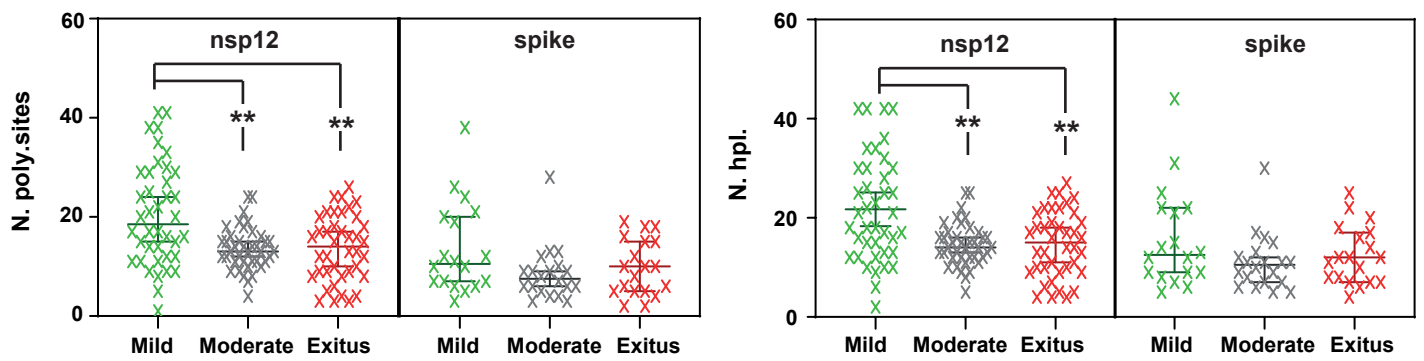




Figure 5

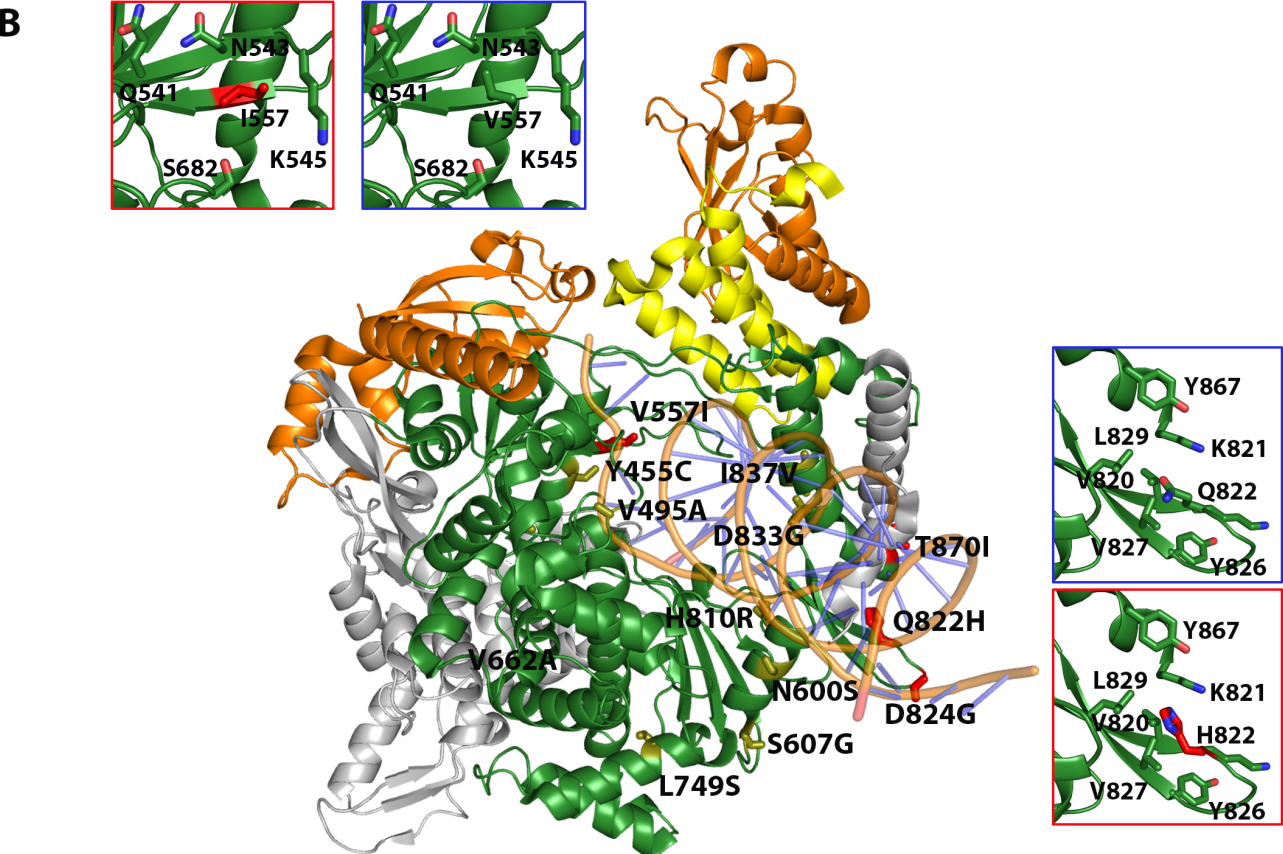
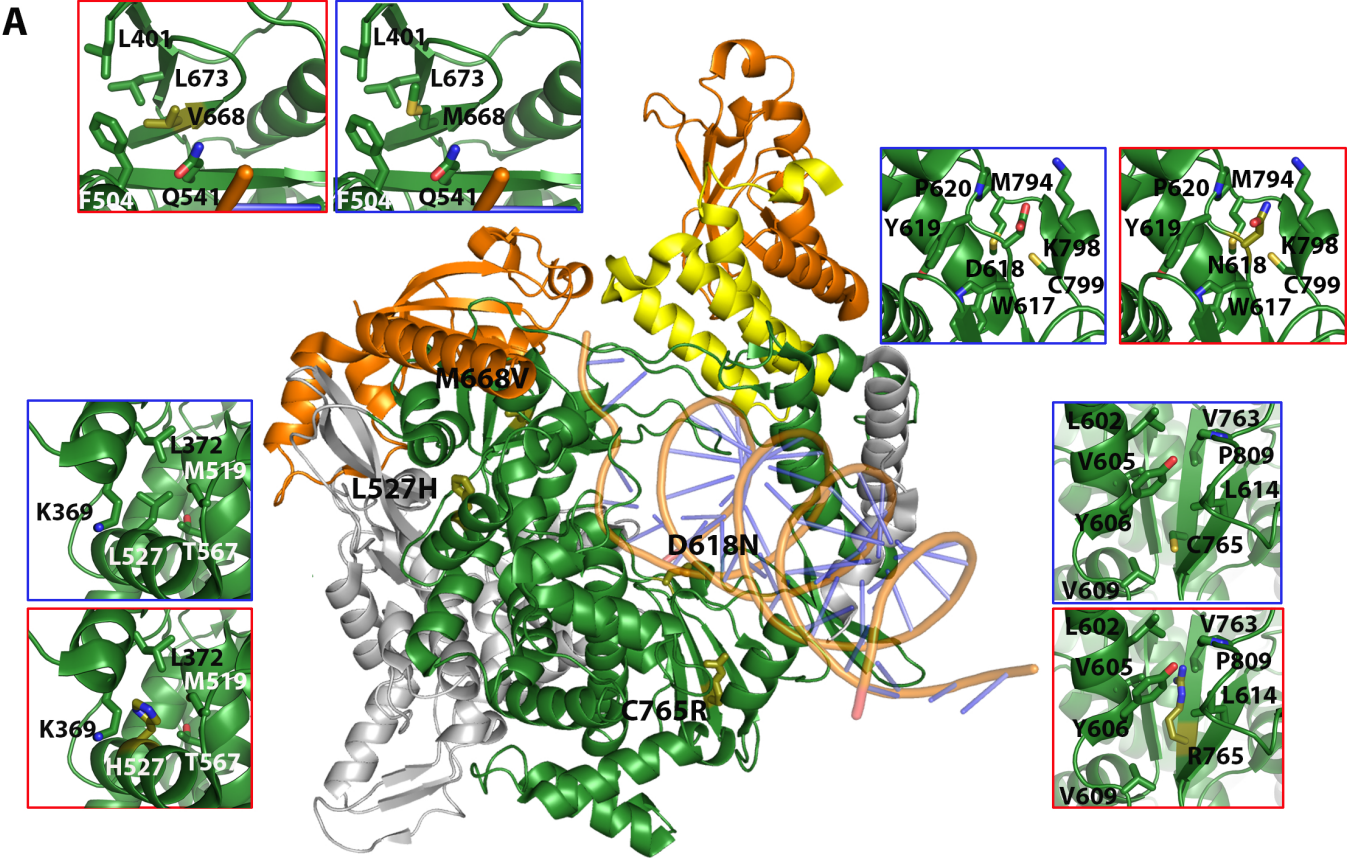


Figure 6

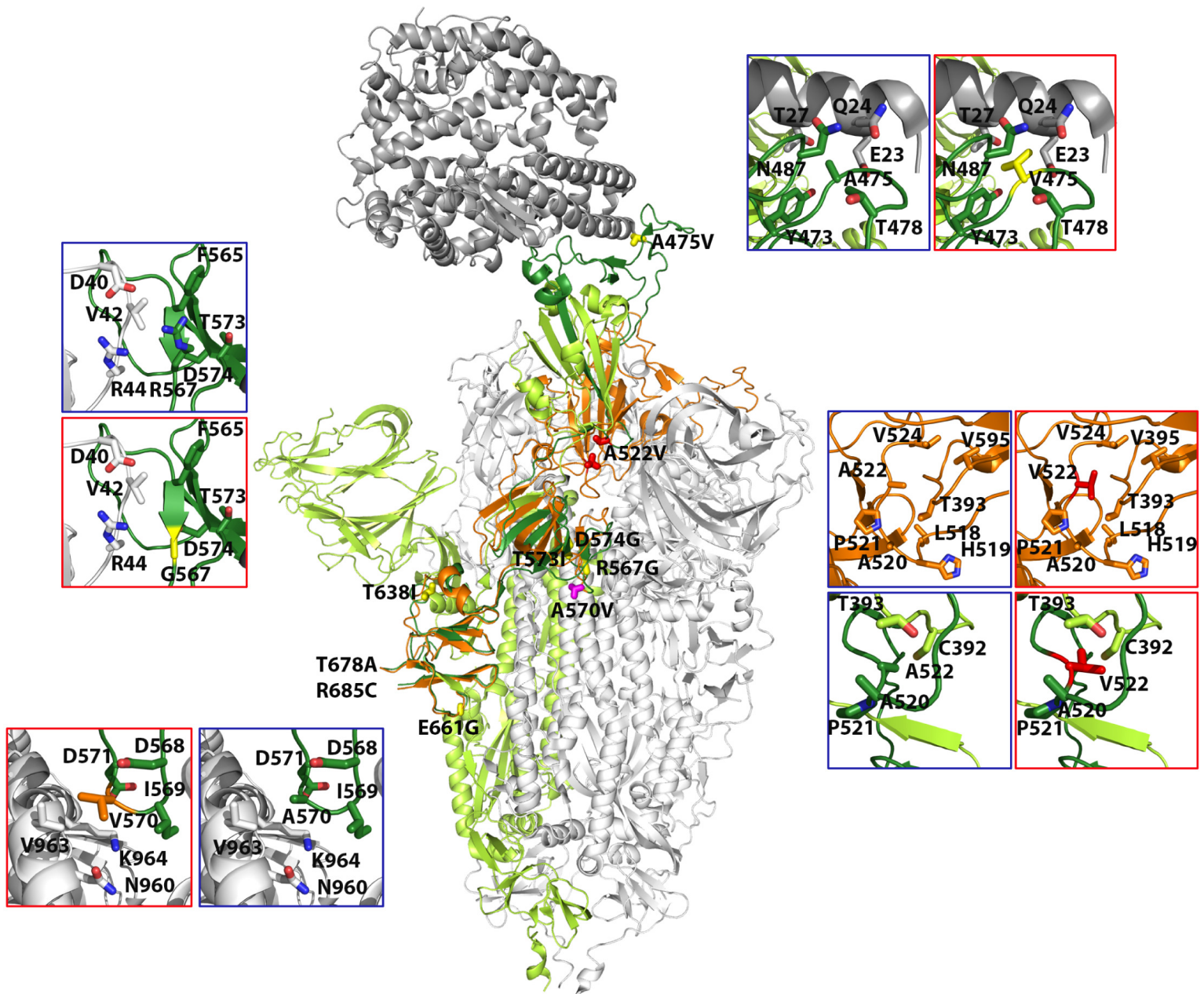
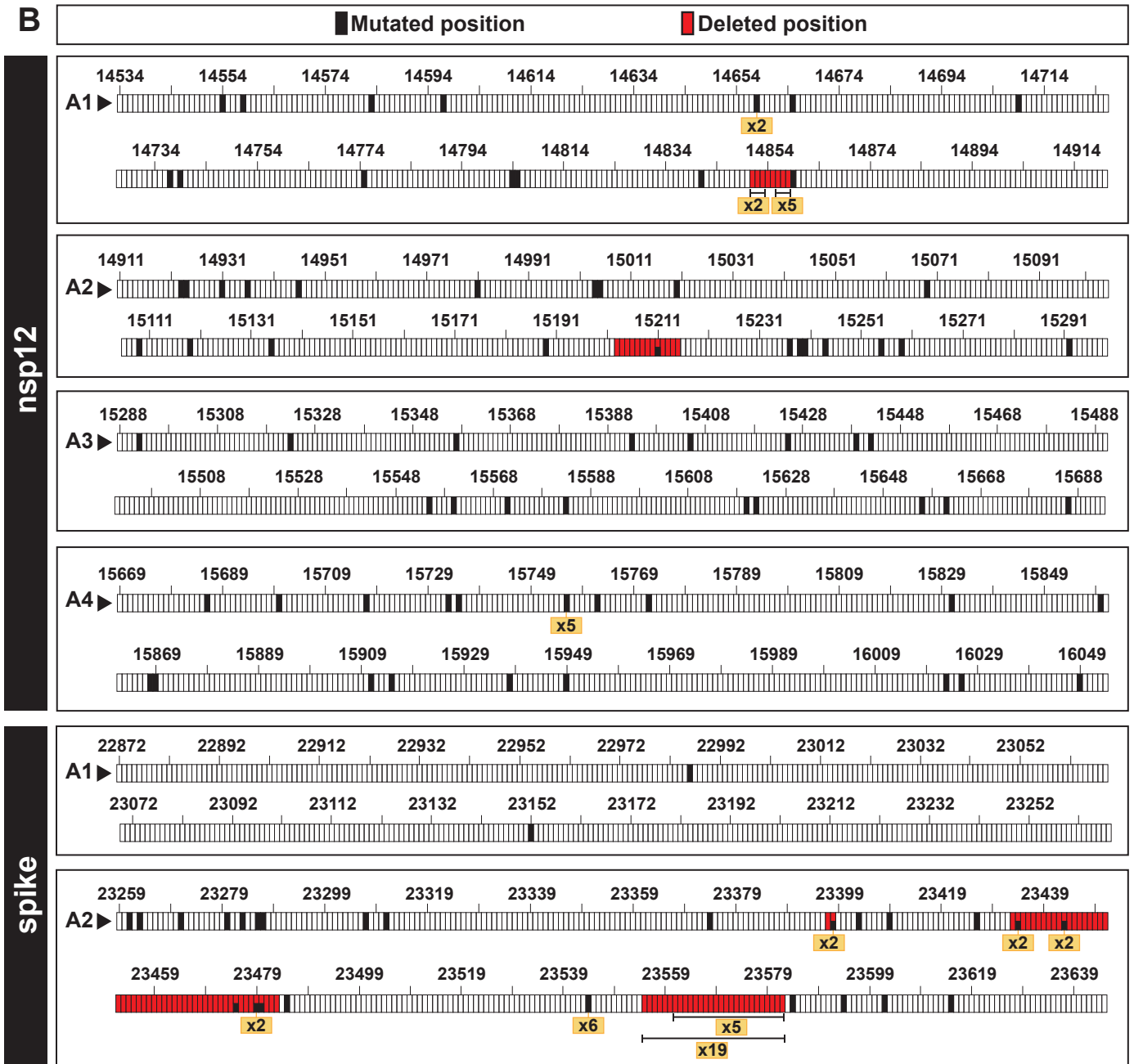
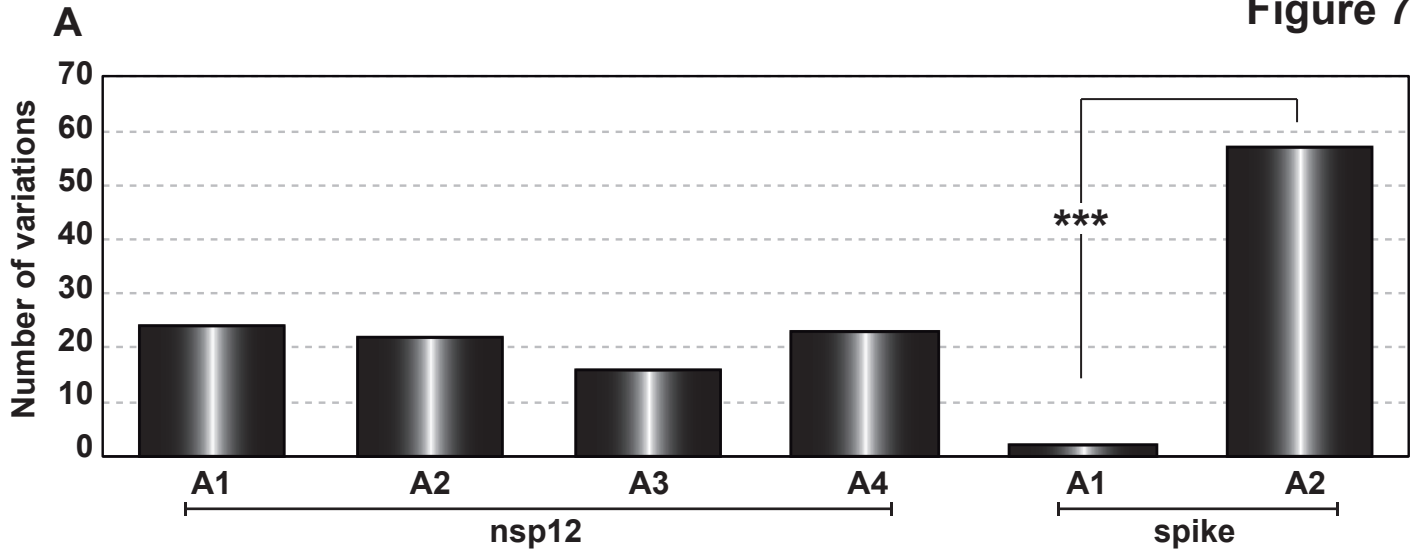


Figure 7



**Table 1.** Point mutations in the mutant spectra of SARS-CoV-2 isolates<sup>a</sup>.

|              |                                 | <b>Disease severity</b> |             |                 |               |
|--------------|---------------------------------|-------------------------|-------------|-----------------|---------------|
|              |                                 | <b>Total</b>            | <b>Mild</b> | <b>Moderate</b> | <b>Exitus</b> |
| <b>nsp12</b> | <b>Transitions (%)</b>          | 68 (97.14%)             | 56 (96.55%) | 3 (100%)        | 10 (100%)     |
|              | <b>Transversions (%)</b>        | 2 (2.86%)               | 2 (3.45%)   | 0 (0%)          | 0 (0%)        |
|              | <b>p-value</b>                  | <0.001                  | <0.001      | 0.051           | <0.001        |
|              | <b>Significance<sup>b</sup></b> | ***                     | ***         | n.s.            | ***           |
|              | <b>Synonymous (%)</b>           | 29 (41.43%)             | 24 (41.38%) | 2 (66.67%)      | 4 (40%)       |
|              | <b>Non-Synonymous (%)</b>       | 41 (58.57%)             | 34 (58.62%) | 1 (33.33%)      | 6 (60%)       |
|              | <b>p-value</b>                  | 0.031                   | 0.047       | 0.5             | 0.327         |
|              | <b>Significance<sup>b</sup></b> | *                       | *           | n.s.            | n.s.          |
| <b>spike</b> | <b>Transitions (%)</b>          | 26 (96.30%)             | 24 (100%)   | 3 (100%)        | 0 (0%)        |
|              | <b>Transversions (%)</b>        | 1 (3.70%)               | 0 (0%)      | 0 (0%)          | 1 (100%)      |
|              | <b>p-value</b>                  | <0.001                  | <0.001      | 0.051           | 0.5           |
|              | <b>Significance<sup>b</sup></b> | ***                     | ***         | n.s.            | n.s.          |
|              | <b>Synonymous (%)</b>           | 11 (40.74%)             | 10 (41.67%) | 1 (33.33%)      | 0 (0%)        |
|              | <b>Non-Synonymous (%)</b>       | 16 (59.26%)             | 14 (58.33%) | 2 (66.67%)      | 1 (100%)      |
|              | <b>p-value</b>                  | 0.138                   | 0.193       | 0.5             | 0.051         |
|              | <b>Significance<sup>b</sup></b> | n.s.                    | n.s.        | n.s.            | n.s.          |

<sup>a</sup> Different number of point mutations distributed according to COVID-19 severity in the nsp12 (polymerase)- and spike-coding region.

<sup>b</sup> The statistical significance of the differences (n.s., not significant; \* p<0.05; \*\*\* p<0.001) was calculated using the proportion test.

**Table 2.** Amino acid substitutions at the nsp12 (polymerase) in the mutant spectra of SARS-CoV-2<sup>a</sup>.

| <b>Low frequency substitutions (0.5% - 2%)</b> |                     |               |                      |   |
|--|---------------------|---------------|----------------------|---|
| <b>Patient category</b>                        | <b>Substitution</b> | <b>PAM250</b> | <b>SNAP2 (score)</b> | <b>Location and possible structural or functional effects</b>   |
| <b>Mild</b>                                    | V373A               | 0             | Neutral (-55)        | Interface, between NiRan and fingers. Loss of a side chain that may interact with L527 and I536 (at 4Å and 3.8Å, respectively), of fingers domain.  |
|  | D499G               | 1             | Effect (70)          | RNA template binding region, but not in directly contact with RNA. May enhance RNA binding through increase in electropositivity.   |
|  | L514P               | -3            | Neutral (-34)        | Near V83 of nsp7. Could affect the interaction between nsp7 and nsp12-polymerase, although other nsp12 residues (F368, L372, F506) are also involved.   |
|  | L527H               | -2            | Effect (51)          | Fingers' helix in contact with the NiRan. It may require structural accommodation in a hydrophobic environment.   |
|  | V560A               | 0             | Effect (3)           | Palm, motif B Side chain of V560 interacts with S681, generating the Up/Down positioning of loop B, involved in RNA translocation, as described in picornaviruses (Sholders et al., 2014). The V-A substitution would inhibit this interaction. |
|  | D618N               | 2             | Effect (75)          | Catalytic D of motif A. Loss of polymerization function.  |
|  | N628S               | 1             | Neutral (-35)        | Fingers. Establishes links with a helix and a loop from fingers through salt bridges. S628 breaks the links and may increase domain flexibility.  |
|  | M668V               | 2             | Neutral (-54)        | Exposed residue in the template entry channel. Substitution M-V would lead to an expansion of the channel.  |
|  | L727P               | -3            | Effect (66)          | Lower part of palm domain. A P residue in this position fits well into a region rich in aromatic amino acids.   |
|  | C765R               | -4            | Effect (76)          | β- strand of the hairpin forming motif A that includes the active site. R at this position would disrupt the surroundings of the active site, probably inducing a non-functional protein.   |
| <b>Exitus</b>                                  | L372F               | 2             | Effect (8)           | Interface, between NiRan and fingers. F may reinforce the hydrophobic environment.  |
| <b>Medium frequency substitutions (2%-30%)</b> |                     |               |                      |   |
| <b>Patient category</b>                        | <b>Substitution</b> | <b>PAM250</b> | <b>SNAP2 (score)</b> | <b>Location and possible structural or functional effects</b>   |

|               |       |   |               |   |
|---------------|-------|---|---------------|---|
| <b>Exitus</b> | V557I | 4 | Neutral (-51) | Close to the entry of the RNA template channel. In contact with the nitrogen base T+1. An I may enhance the stability of this connection. |
|---------------|-------|---|---------------|---|

### High frequency substitutions (>90%)

| <b>Patient category</b> | <b>Substitution</b> | <b>PAM250</b> | <b>SNAP2 (score)</b> | <b>Location and possible structural or functional effects</b>     |
|-------------------------|---------------------|---------------|----------------------|---|
| <b>Exitus</b>           | Q822H               | 3             | Neutral (-85)        | In a loop of the thumb domain. An H could enhance loop stability. |

<sup>a</sup>The sequenced region spans amino acids 366 to 871. Substitutions are divided according to the frequency at which they are found in the mutant spectra, and disease category [mild, moderate or severe (exitus) as defined in Methods] (Figure 2). PAM250 and SNAP2 scores have been calculated as described in [49] and [50], respectively. Possible structural effects have been predicted from the location of the substitution in the three-dimensional structure of nsp12 (polymerase) (Figure 5).

**Table 3.** Amino acid substitutions at the spike (S) protein in the mutant spectra of SARS-CoV-2<sup>a</sup>.

| <b>Low frequency substitutions (0.5% - 2%)</b> |                     |               |                      |  |
|--|---------------------|---------------|----------------------|--|
| <b>Patient category</b>                        | <b>Substitution</b> | <b>PAM250</b> | <b>SNAP2 (score)</b> | <b>Location and possible structural or functional effects</b>  |
| <b>Mild</b>                                    | R567G               | -3            | Effect (47)          | Contact region, involved in the formation of the S trimers. This substitution would eliminate the R567-D40 salt bridge, involved in regulation of the viral fusion. The R-G substitution could facilitate fusion with cell.  |
|  | T573I               | 0             | Neutral (-54)        | $\beta$ -chain next to R567 and close to a V and two F residues. The T-I substitution may strengthen hydrophobic contacts in the region.   |
|  | D574G               | 1             | Neutral (-39)        | $\beta$ -chain next to T573. Loss of contact with K557 and increase of flexibility.  |
| <b>Mild and Moderate</b>                       | E661G               | 0             | Effect (41)          | Exposed residue that could interact with Q779 of another chain in the S trimer. A G would prevent this interaction.  |
| <b>Medium frequency substitutions (2%-30%)</b> |                     |               |                      |  |
| <b>Patient category</b>                        | <b>Substitution</b> | <b>PAM250</b> | <b>SNAP2 (score)</b> | <b>Location and possible structural or functional effects</b>  |
| <b>Mild</b>                                    | A475V               | 0             | Neutral (-88)        | Interaction with ACE2. V may increase contact with ACE2 receptor.  |
|  | T678A               | 1             | Neutral (-23)        | Loop near the furin cleavage site. Expected to be exposed upon furin cleavage. However, this region appears disordered in the deposited structures.  |
|  | R685C               | -4            | Effect (26)          | Furin cleavage site (PRRAR). A C would either inhibit the cleavage or decrease the efficacy of the excision, thus hindering the S1/S2 excision.  |
| <b>Moderate</b>                                | A570V               | 0             | Neutral (-88)        | Interaction region to form S trimers. V could bring closer the two chains due to its larger and more hydrophobic side-chain.   |
| <b>High frequency substitutions (&gt;90%)</b>  |                     |               |                      |  |
| <b>Patient category</b>                        | <b>Substitution</b> | <b>PAM250</b> | <b>SNAP2 (score)</b> | <b>Location and possible structural or functional effects</b>  |
| <b>Exitus</b>                                  | A522V               | 0             | Neutral (-71)        | Loop close to the hinge, linking the RBD and the sub-domain 1 of S1. This loop facilitates the transition from the “open” to the “erect” position of the RBD. The A-V substitution may enhance the stability of the RBD open position, due to its proximity to other hydrophobic residues. |

<sup>a</sup>The sequenced region spans amino acids 438 to 694. Substitutions are divided according to the frequency at which they are found in the mutant spectra, and disease category [mild, moderate or severe (exitus) as defined in Methods] (Figure 2). PAM250 and SNAP2 scores have been calculated as described in [49] and [50], respectively. Possible structural effects have been predicted from the location of the substitution in the three-dimensional structure of S (Figure 6).