

1 **Multi-Omics and Integrated Network Approach to Unveil Evolutionary Patterns,**  
2 **Mutational Hotspots, Functional Crosstalk and Regulatory Interactions in SARS-CoV-2**

3

4 Vipin Gupta<sup>¥1</sup>, Shaiza Haider<sup>¥2</sup>, Mansi Verma<sup>¥3</sup>, Kalaiarasan Ponnusamy<sup>4</sup>, Md. Zubair  
5 Malik<sup>5</sup>, Nirjara Singhvi<sup>6</sup>, Helianthous Verma<sup>7</sup>, Roshan Kumar<sup>8</sup>, Utkarsh Sood<sup>9</sup>, Princy Hira<sup>10</sup>,  
6 Shiva Satija<sup>3</sup> and Rup Lal<sup>9\*</sup>

7

8 <sup>1</sup> PhiXGen Private Limited, Gurugram, Haryana 122001, India

9 <sup>2</sup> Department of Biotechnology, Jaypee Institute of Information Technology, Noida, sector-62,  
10 Uttar Pradesh, India

11 <sup>3</sup> Department of Zoology, Sri Venkateswara College, University of Delhi, New Delhi-110021,  
12 India

13 <sup>4</sup> School of Biotechnology, Jawaharlal Nehru University, New Delhi, India.

14 <sup>5</sup> School of Computational and Integrative Sciences, Jawaharlal Nehru University, New Delhi,  
15 India.

16 <sup>6</sup> Department of Zoology, University of Delhi, New Delhi-110007, India

17 <sup>7</sup> Molecular Biology and Genomics Research Laboratory, Ramjas College, University of  
18 Delhi, 110007, India

19 <sup>8</sup> P.G. Department of Zoology, Magadh University, Bodh Gaya, Bihar-824234, India

20 <sup>9</sup> The Energy and Resources Institute, Darbari Seth Block, IHC Complex, Lodhi Road, New  
21 Delhi-110003, India

22 <sup>10</sup> Maitreyi College, University of Delhi. Chanakyapuri, New Delhi 110021

23

24 \*Corresponding Author

25 <sup>¥</sup> Contributed Equally

26

27 Corresponding author Email: [ruplal@gmail.com](mailto:ruplal@gmail.com)

28

29

30

31

32

### 33 **Abstract**

34 SARS-CoV-2 responsible for the pandemic of the Severe Acute Respiratory Syndrome  
35 resulting in infections and death of millions worldwide with maximum cases and mortality in  
36 USA. The current study focuses on understanding the population specific variations  
37 attributing its high rate of infections in specific geographical regions which may help in  
38 developing appropriate treatment strategies for COVID-19 pandemic. Rigorous phylogenetic  
39 network analysis of 245 complete SARS-CoV-2 genomes inferred five central clades named a  
40 (ancestral), b, c, d and e (subtype e1 & e2) showing both divergent and linear evolution types.  
41 The clade d & e2 were found exclusively comprising of USA strains with highest known  
42 mutations. Clades were distinguished by ten co-mutational combinations in proteins; Nsp3,  
43 ORF8, Nsp13, S, Nsp12, Nsp2 and Nsp6 generated by Amino Acid Variations (AAV). Our  
44 analysis revealed that only 67.46 % of SNP mutations were carried by amino acid at  
45 phenotypic level. T1103P mutation in Nsp3 was predicted to increase the protein stability in  
46 238 strains except six strains which were marked as ancestral type; whereas com (P5731L &  
47 Y5768C) in Nsp13 were found in 64 genomes of USA highlighting its 100% co-  
48 occurrence. Docking study highlighted mutation (D7611G) caused reduction in binding of  
49 Spike proteins with ACE2, but it also showed better interaction with TMPRSS2 receptor  
50 which may contribute to its high transmissibility in USA strains. In addition, we found host  
51 proteins, MYO5A, MYO5B & MYO5C had maximum interaction with viral hub proteins  
52 (Nucleocapsid, Spike & Membrane). Thus, blocking the internalization pathway by inhibiting  
53 MYO-5 proteins which could be an effective target for COVID-19 treatment. The functional  
54 annotations of the Host-Pathogen Interaction (HPI) network were found to be highly  
55 associated with hypoxia and thrombotic conditions confirming the vulnerability and severity  
56 of infection in the patients. We also considered the presence of CpG islands in Nsp1 and N  
57 proteins which may confers the ability of SARS-CoV-2 to enter and trigger methyltransferase  
58 activity inside host cell.

### 59 **Introduction**

60 In December 2019, a novel RNA virus, Severe Acute Respiratory Syndrome Corona Virus-2  
61 (SARS-CoV-2), belonging to Coronaviridae family (betacoronavirus), emerged as the reason  
62 for the chaos of pneumonia disease also called Covid-19 in Chinese city, Wuhan (Li et al.,  
63 2020). Covid-19 was declared as pandemic by WHO on March 11, 2020 (Astuti and Ysrafil,  
64 2020). Major outbreaks were reported in many locations of China, USA, Italy, Spain, Japan,

65 and South Korea. As of date it has already spread to more than 200 countries of the world  
66 surpassing more than 30 thousand deaths and 6 million reported active cases worldwide  
67 (<https://www.worldometers.info/coronavirus/>).

68 SARS-CoV-2 is a single stranded RNA virus with a genome size ranging from 29.8 kb  
69 to 29.9 kb (Khailany et al., 2020). The genomic repertoire of SARS-CoV-2 comprises of 10  
70 open reading frames (ORFs) encoding 27 proteins (Abduljalil and Abduljalil, 2020). ORF1ab  
71 encodes for 16 non-structural proteins (Nsp) whereas structural proteins include spike (S),  
72 envelope (E), membrane (M), and nucleocapsid (N) proteins (Pyrce et al., 2007; Yang and  
73 Leibowitz, 2015). In addition, the genome of SARS-CoV-2 comprises of ORF3a, ORF6,  
74 ORF7a, ORF7b, ORF8 and ORF9 genes encoding six accessory proteins, flanked by 5' and 3'  
75 UTRs (Khailany et al., 2020). In our previous study (Kumar et al., 2020), a higher mutational  
76 rate in the genomes from different geographical locations around the world by accumulation  
77 of Single Nucleotide Polymorphisms (SNP) was reported. Even during these early stages of  
78 the global pandemic, genomic surveillance has been used to differentiate circulating strains  
79 into distinct, geographically based lineages (Forster et al., 2020). However, the ongoing  
80 analysis of this global dataset suggests no consolidated significant links between SARS-CoV-  
81 2 genome sequence variability, virus transmissibility and disease severity.

82 It is known that mutations at both genomic and protein level are “Hormonical Orchestra” (Yu  
83 et al., 2019) that drives the evolutionary changes, demanding a detailed study of SARS-CoV-2  
84 mutations to understand its successful invasion and infection. The study analyzed that  
85 mutational profiles of SARS-CoV-2 isolates show very high mutational rates that show the  
86 isolates more virulent, causing significant harm to the hosts (Mandal et al., 2020). Thus, in the  
87 present study, we selected 245 genomic sequences of SARS-CoV-2 deciphering the  
88 phylogenetic relationships, tracing them to SNPs at nucleotide and amino acid (Amino Acid  
89 Variation) levels and performing structural re-modelling. Our results revealed the  
90 evolutionary relationships among the strains predicting Nsp3 as mutational hotspot for SARS-  
91 CoV-2. We further extended the study to understand mechanism of host immunity evasion by  
92 Host-Pathogen Interaction (HPI) and confirming their interactions with host proteins by  
93 docking studies. We identified sparsely distributed hubs which may interfere and control  
94 network stability as well as other communities/modules. This indicated the affinity to attract a  
95 large number of low-degree nodes toward each hub, which is a strong evidence of controlling  
96 the topological properties of the network by these few hubs (Nafis et al., 2015). We also  
97 analyzed the transfer of genomic SNPs to amino acid levels and associations of CpG islands

98 contributing towards the pathogenicity of SARS-CoV-2. The existence of CpG islands has  
99 always been connected with the epigenetic regulation and act as hotspots for methylation  
100 (Jones, 2012; Shiraishi et al., 2002; Hoelzer et al., 2008). Here also, the conservancy found in  
101 possession of CpG islands towards the extremities of all the genomes considered in the  
102 present analysis indicate their importance in evading host immunity. Our study showed an  
103 overall depiction of SARS-CoV-2 variations and interactions that eventually may lead to  
104 development of rational therapeutic measures and medication against COVID-19.

## 105 **Material and Methods**

### 106 **Selection of genomes, annotations and phylogeny construction**

107 Publicly available genomes of SARS-CoV-2 viruses were obtained from the NCBI database  
108 (<https://www.ncbi.nlm.nih.gov/genbank/sars-cov-2-seqs/>). Until March 31, 2020 only 375  
109 SARS-CoV-2 genomes were available in the databases. The data was screened for unwanted  
110 ambiguous bases using N-analysis program, based on which 245 complete and clean genomes  
111 of SARS-CoV-2 were selected for further analysis (Supplementary Info 1). A manually  
112 annotated reference database was generated using GenBank file of severe acute respiratory  
113 syndrome coronavirus 2 isolate- SARS-CoV-2/SH01/human/2020/CHN (Accession number:  
114 MT121215.1) and open reading frames (ORFs) were predicted against the formatted database  
115 using prokka (-gcode 1) (Seemann, 2014). Genomic sequences included in the analysis  
116 belongs to different countries namely, USA (168), China (53), Pakistan (2), Australia (1),  
117 Brazil (1), Finland (1), India (2), Israel (2), Japan (5), Vietnam (2), Nepal (1), Peru (1), South  
118 Korea (1), Spain (1), Sweden (1). Whole genomes nucleotide and protein sequences were  
119 aligned using mafft (Katoh et al., 2013) at 1000 iterations. The alignments so obtained were  
120 processed for phylogeny construction using BioEdit software (Hall et al., 2011). The  
121 nucleotide-based phylogeny was annotated and visualized on iTOL server (Letunic and Bork,  
122 2006). While, amino acid-based phylogeny was visualized and annotated using GrapeTree  
123 (Zhou et al., 2018).

124

### 125 **Genotyping based on SNP/AAV**

126 To detect nucleotide and amino acid variations (AAV) among 245 genomes of SARS-CoV-2,  
127 sequence alignment of nucleotide and amino acid, respectively were performed against the  
128 reference genome. The change of nucleotide and amino acid was calculated as point variations  
129 and were recorded. The interpolation and visualization were plotted using computer programs

130 in Python. Co-mutation were predicted and clustering was performed using MicroReact  
131 (Argimon et al., 2016)

132

### 133 **Data and Computer programs:**

134 The genomic analytics is performed using programs in Python and Biopython libraries (Cock  
135 et al., 2009). The computer programs and the updated SNP profiles of SARS-CoV-2 isolates  
136 are available upon requests.

137

### 138 **Construction of the Host-Pathogen Interaction Network of SARS-CoV-2**

139 In order to find the HPI, we subjected SARS-CoV-2 proteins to Host-Pathogen interaction  
140 databases such as Viruses.STRING v10.5 (Cook et al., 2018) and HPIDB3.0 (Ammari et al.,  
141 2016) for predicting their direct interaction with human as the principal host. The HPI  
142 network was constructed and visualized using Cytoscape v3.7.2 (Shannon, et al., 2003). In the  
143 constructed Network, proteins with highest degree, which interact with several other signaling  
144 proteins in the network indicate a key regulatory role as a hub. In our study, using  
145 NetworkAnalyser (Assenov, et al., 2008), plugin of Cytoscape v3.7.2, we identified the hub  
146 protein and subjected to functional analysis. The network was functionally annotated using  
147 STRINGApp and StringEnrichment app (Doncheva et al., 2019) plugin of Cytoscape using  
148 Reactome, GO, InterPRO, KEGG and Pfam databases. This analysis provides an opportunity  
149 of a more precise understanding of the biological functions, providing valuable clues for  
150 biologists.

151

### 152 **Computational structural analysis on wild-type and mutant SARS-CoV-2 proteins**

153 SARS-CoV-2 proteins sequences were retrieved from the NCBI genome database and  
154 pairwise sequence alignment of wild-type and mutant proteins were carried out by the Clustal  
155 Omega tool (Sievers et al., 2011). The wild-type and mutant homology model of S-protein,  
156 NspNsp12 and Nsp13 were constructed using the SWISSMODEL (Waterhouse et al., 2018),  
157 whereas the 3D structure of ORF8, ORF3A, Nsp2, Nsp3 and Nsp6 were predicted using  
158 Phyre2 server (Kelley et al., 2015). The host proteins (TMPRSS2, RPS6, ATP6V1G1 and  
159 MYO5C) 3D structures were generated using the SWISSMODEL and ACE2 structure  
160 retrieved from the PDB database (PDB ID: 6M17). These structures were energy minimized  
161 by the Chiron energy minimization server (Ramachandran et al., 2011). The effect of the  
162 mutation was analyzed using HOPE (Venselaar et al., 2010) and I-mutant (Capriotti et al.,  
163 2006). The I-mutant method allows us to predict the stability of the protein due to mutation.

164 The docking studies for wild and mutant SARS-CoV-2 proteins with host proteins was carried  
165 out using PatchDock Server (Schneidman-Duhovny et al., 2005). Structural visualizations and  
166 analysis were carried out using pyMOL2.3.5 (Jacobson et al., 2002).

### 167 **Analysis of CpG regions**

168 SARS-CoV-2 genomes were analysed for the presence of CpG regions that can be targeted for  
169 methylation induced gene silencing. To locate the CpG regions, meth primer 2.0  
170 (<http://www.urogene.org/methprimer2/>) and the CpG Plot  
171 (<http://www.ebi.ac.uk/Tools/emboss/cpgplot/>) programs were used, although some variations  
172 were found in both the programs. Both the programs were run on default parameters of a  
173 sequence window longer than 100 bp; GC content of  $\geq 50\%$ , and an observed/expected CpG  
174 dinucleotide ratio  $\geq 0.60$ . The presence of common CpG islands was confirmed by performing  
175 BLAST using the above reference strain.

176

## 177 **Results and Discussion**

### 178 **Phylogenetic relationship between different SARS-CoV-2 strains**

179 In our previous study, we reported a mosaic pattern of phylogenetic clustering of 95 genomes  
180 of SARS-COV-2 isolated from different geographical locations (Kumar et al., 2020). Strains  
181 belonging to one country were found clustered with distant countries strains but not with the  
182 neighboring one. Taking clue from these studies we constructed phylogenetic relatedness of  
183 245 strains of SARS-COV-2 from USA, China, and several other countries including, Spain,  
184 Vietnam, Peru, Finland and Pakistan and unravel the significant association of evolutionary  
185 patterns among SARS-CoV-2 based on their geographical locations predicting their mosaic  
186 phylogenetic arrangements. It was found that the majority of strains from USA were clustered  
187 together, but comparatively high divergences were found in strains isolated from China and  
188 Japan. Japanese strains were found to be scattered and formed clusters with strains from USA,  
189 Pakistan, Vietnam, Taiwan, and China. Even with less number of genomes sequences from  
190 Japan, Vietnam and Peru revealed a highly scattered pattern and formed close associations  
191 with that of USA and Chinese strains. Strains reported from patients of Taiwan (MT192759),  
192 Australia (MT007544), South Korea (MT039890), Nepal (MT072688) and Vietnam  
193 (MT192773, MT192772) had travel histories from Wuhan, China (Cheng et al., 2020).  
194 However, a strain from Pakistan (MT240479) which clustered with the Japanese strains was  
195 found to be isolated from patient having travel history from Iran. Indian strains (MT050439,  
196 MT012098) that were isolated from patients who travelled from Dubai, clustered with

197 Chinese strains. Later, reports confirmed many cases of SARS-CoV-2 in Dubai from China  
198 ([https://www.newsbytesapp.com/timeline/India/58169/271167/coronavirus-2-positive-cases-](https://www.newsbytesapp.com/timeline/India/58169/271167/coronavirus-2-positive-cases-detected-in-delhi-telangana)  
199 [detected-in-delhi-telangana](https://www.newsbytesapp.com/timeline/India/58169/271167/coronavirus-2-positive-cases-detected-in-delhi-telangana)). Thus, a clear landscape of phylogenetic relationships could be  
200 obtained reflecting mosaic clustering patterns in accordance with the travel history of patients  
201 (Figure1A). However, results were in contradiction with the genomic analysis of SARS-  
202 COV-2 by Foster et al.,( 2020) where they predicted the linear/directive evolution from  
203 ancestral node a to node b and c. Whereas we report here both divergent (from ancestral node  
204 a to b, c & e) and directive (node c to d) evolution among the SARS-CoV-2 strains (Figure1B  
205 and Figure 3B). Since genome-based phylogeny did not highlight the amino acid level  
206 changes, thus to ascertain the variations among the SARS-CoV-2 strains at phenotypic level,  
207 we constructed whole proteome alignment-based phylogeny, clustered the 245 strains into  
208 five major **clades a-e** (Figure 1B). The first cluster, Clade-a had maximum nodes (46),  
209 including reference node, and strains from Nepal (MT072688), Pakistan (MT262993), Taiwan  
210 (MT192759) along with 15 strains from USA and 27 strains from China. It also had the  
211 mutated daughter nodes (highlighted by # in figure 3 B for corresponding nodes) radiating  
212 outwards, belonging to China, Finland (MT020781), India (MT012098), Japan (LC534419,  
213 LC529905), Taiwan (MT066176), Vietnam (MT192772-3), Brazil (MT126808), Australia  
214 (MT007544), South Korea (MT039890) and Sweden (MT093571) along with seven USA  
215 strains (Figure 1B). This clade represented the ancestral node as it harbored the oldest known  
216 SARS-CoV-2 strain from China and laid down the foundation of rest of the mutated daughter  
217 strains worldwide, marking the onset of the divergence in SARS CoV-2. Three significantly  
218 diverged network nodes originated from the ancestral clade-a and were marked as clade-b, c  
219 and e (Figure 1B). For **Clade-b**, central node included only four strains in which two were  
220 from USA (MT184912, MT276328) and one each from Israel (MT276597) and Japan  
221 (LC528233). Its major descended radiant belonged to Japan (LC528232, LC534418), Pakistan  
222 (MT240479), USA (MT184913, MT184910, MN997409) and China (MT049951,  
223 MT226610). It was observed that one of the Chinese strains in clad-b (MT226610) had the  
224 longest branch length making the strain very distinct (harboring 25 mutations) by showing  
225 exceptionally high rate of evolution. In **Clade-c** lineage, small central node was comprised of  
226 Taiwan (MT066175), USA (MT246667, MT233526, MT020881, MT985325, MT020880)  
227 and Chinese (MN938384, LR757995) strains. Interestingly one strain each from Spain  
228 (MT233523) and India (MT050493) were also found radiating as daughter node from the  
229 central one. **Clade-d** lineage, which was originated from clade-c lineage, consisted only of  
230 USA strains both in central nodes and radiations. Importantly, 2 strains (MT263416,

231 MT246471) were found most divergent with varied mutation suggesting the high rate of  
232 evolution among USA strains which might be linked with the high pathogenicity among them.  
233 **Clade-e** bifurcated into two sub-clads (e1 and e2) by significant set of mutations. Sub-clad-e1  
234 include six strains from USA, one from Israel (MT276598) with radiating nodes from Peru  
235 (MT263074) and USA (MT276327); whereas, sub-clad e2 had 32 strains belonging to USA.  
236 Thus, formation of five major evolutionary clades and subclades based on the amino acid  
237 phylogeny needs attention for identifying the assessment of divergence among SARS-CoV-2  
238 strains. This divergence is a proof of the random evolution of SARS-CoV-2 suggesting  
239 network expansion in five clads contradicting to the earlier directed evolution proposed by  
240 Foster et al., 2020.

241

### 242 **Genotyping and variation estimation**

243 In order to understand the implication of mosaic pattern of transmissions and evolutionary  
244 lineage clustering (Clad a-e), we studied the Single Nucleotide Polymorphism (SNP)  
245 genotyping from the 245 genome sequences as mutation counts along with their frequency at  
246 specific genomic locations. Mutational changes at phenotypic levels were also weighed by  
247 assessing Amino Acid Variations (AAV). Interpolations of the SNPs/AAVs data were made  
248 by assessing their frequency, genomic positions and type of SNPs/AAVs (Figure 2B),  
249 highlighted a large mutational diversity among the virus isolates. We identified a total of 12  
250 SNP types (A>G, A>C, A>T, C>A, C>G, C>T, G>A, G>C, G>T, T>A, T>C, T>G) accounting  
251 for mutations at 297 genomic locations (Figure 2A, 2B). Overall pattern of SNPs suggested  
252 C>T transition as the most common mutation in the entire genomic sets (Figure 2A), however  
253 highest frequency was recorded for T>C transitions (Figure 2B). Based on the genomic  
254 arbitrators SNP frequencies, we analyzed 14 major locations inside the genomes of SARS-  
255 CoV-2 for potential mutation generating different allelic forms for genes (Table 1). The SNP  
256 of C>T was first observed at 67<sup>th</sup> location in 5' UTR region of leader sequence with a  
257 frequency of 45 followed by Nsp2 at two locations (885 & 2863) with the frequency of 29 and  
258 44, respectively. Nsp3/PL-PRO and Nsp8 marked the highest frequency of 238 SNP counts of  
259 T>C at 5852 and 12299 locations. Another T>C SNP was observed in ORF8 with frequency  
260 of 88 at 27973 location. C>T SNP transformation was found in Nsp4 and Nsp12 with the  
261 frequency of 88 and 44 at location 8608 and 14234, respectively. Non-structural protein,  
262 Nsp13 was strangely found harboring two different SNP (C>T and A>G) at three different  
263 locations (17573, 17684, 17886) with a relatively high frequency of 68, 63 and 63



264 respectively. A>G SNP conversion in S (Spike) protein was found with a frequency of 43. A  
265 Low SNP count of G>T transitions were falling in the ORF3a and Nsp6 with frequency of 32  
266 and 21, respectively (Table 1). Though, all SNP counts do not reflect the phenotypic change at  
267 protein level and therefore must be estimated at the translation levels for their significant  
268 phenotypic effect. Although 297 genomic locations harbored SNPs but their corresponding  
269 AAV were found only in 200 genomic locations accounting for 67.34% conversion efficiency.  
270 Out of 14 high frequency SNPs, only 9 mutations [Nsp2 (T265I), Nsp3 (S1920P), Nsp6  
271 (L3605F), Nsp12 (P4618L), Nsp13 (P5731L, Y5768C), S (D7611G), Orf3a (Q8327H), Orf8  
272 (L9033S)] were found to reflect at protein level with the highest frequency of 238 in Nsp3  
273 (Table 1). These proteins are known to play various regulatory roles and therefore, mutations  
274 at amino acid level can modulate their catalytic activity drastically. Specifically, Nsp3 is the  
275 largest and essential component of replication complex in the SARS-CoV-2 genome (Lei et  
276 al., 2018) and along with Nsp2 it forms a transcriptional complex in endosome of the infected  
277 host cell (Wu et al., 2020). Nsp6 is a multiple-spanning transmembrane protein located into  
278 the ER where they induce autophagosomes via an omegasome intermediate (Cottam et al.,  
279 2014). Interestingly, the mutation of L3605F causes stiffness in the secondary structure of  
280 Nsp6 and leads to low stability of the protein structure in most recent sequences from Asia,  
281 America, Oceania and Europe (Benvenuto et al., 2020). Nsp12 and Nsp13 are the key  
282 replicative enzymes, which require Nsp6, Nsp7 and Nsp10 as cofactors. Nsp12, RNA  
283 dependent RNA polymerase (RdRp) with the presence of the bulkier leucine side chain at  
284 location 4618 is likely to create a greater stringency for base pairing to the templating  
285 nucleotide, thus modulating polymerase fidelity (Sexton et al., 2016). Nsp13 contains a  
286 helicase domain, allowing efficient strand separation of extended regions of double-stranded  
287 RNA and DNA (Pachetti et al., 2020). Dual mutations in Nsp13 were reported with profound  
288 effect on its activity. P5731L, mutation leads to increased affinity of helicase RNA  
289 interaction, whereas Y5768C is a destabilizing mutation increasing the molecular flexibility  
290 and leading to decreased affinity of helicase binding with RNA (Begum et al., 2020).  
291 Therefore, both the mutations were antagonistic in nature. Thus, ORF1ab polyprotein of  
292 SARS-CoV-2 encompasses mutational spectra where signature mutations for Nsp2, Nsp3,  
293 Nsp6, Nsp12 and Nsp13 have been predicted. Amino acid mutations in structural proteins S,  
294 ORF3a and ORF8 have also been observed with varied frequency of 45, 34 and 89  
295 respectively. The mutation in Spike protein (D7611G) has been reported to outcompete other  
296 preexisting subtypes, including the ancestral one. This mutation generates an additional serine  
297 protease (Elastase) cleavage site in Spike protein (Bhattacharya et al., 2020) which is

298 discussed in more details in later sections. ORF3a mutation (Q8327H), is located near TNF  
299 receptor associated factor-3 (TRAF-3) regions and has been reported as molecular differences  
300 marker in many genomes including Indian SARS- CoV-2 genomes (Hassan et al., 2020) for  
301 their delineation. Amino acid change in ORF8 sequence (L9033S) propose that it is preserved  
302 (Koyama et al., 2020) therefore it is critical to examine its biological function in SARS-CoV-  
303 2.

304 Our results showed that the mutations (SNPs and AAV) in the virus were not uniformly  
305 distributed. Genotyping study annotated few mutations in the SARS-CoV-2 genomes at  
306 certain specific locations with high frequency predicting their high selective pressure. Thus,  
307 mutations can be predicted as location-specific but not type-specific by SNP count. Highly  
308 frequent AAV might be associated with the changes in transmissibility and virulence behavior  
309 of the SARS-CoV-2. Therefore, high-frequency AAV mutations in Spike protein, RdRp,  
310 helicase and ORF3a are important factors to consider while developing vaccines against the  
311 fast-evolving strains of SARS-CoV-2.

### 312 **Prevalence of Co-mutation in SARS-COV-2 evolution**

313 Interestingly, we observed co-mutations in Nsp13 at locations 5768 (Nsp13\_1) and  
314 5731(Nsp13\_2) that were prevalent in common 64 genomes, all belonging to USA. The AAV  
315 reported above (Table 1) were further analyzed and found occurring in 10 different  
316 permutations varying from single to multiple mutated protein combinations. Complete details  
317 of these co-mutations combinations are given in Table 2. These co-mutations were mapped  
318 over the divergent phylogeny for indicating the evolutionary divergence among the 245  
319 strains. The phylogram (Figure 1B) showed clear divergence of strains from the parent strain  
320 due to accumulation of mutations at different level of human to human transmission. We  
321 found co-mutations in Nsp3, ORF8, Nsp13, S, Nsp12, Nsp2 and Nsp6 were responsible for  
322 the above divergence.

323 These co-mutations were found linked with lineage **clades a to e**, highlighting their  
324 prevalence among them (Figure.1B). In **clade-a**, 40 genomes harbored mutations at only  
325 Nsp3 protein while six isolates belonging to USA (MT262993, MT044258, MT159716,  
326 MT259248, MT259267) and Pakistan (MT263424) showed no mutation confirming their  
327 lineage same as that of the reference/ancestral genome from China. Therefore, Nsp3 marked  
328 as first mutational hotspot for accumulating amino acid mutations in SARS-CoV-2. Brazil  
329 (MT126808) and USA (MT276331) strains form the descendent from clade-a harbored  
330 Nsp3/Nsp6 as first co-mutation. The **clade-b** also had an additional mutation of ORF8 along

331 with Nsp3 and Nsp6 with three descendant strain from US and China. We observed most  
332 distant Chinese strain (MT226610), clustered in **clade-b** and harbored additional 25 AAV  
333 making it the highly pathogenic strain in the network as reported above in Figure 1B. The  
334 **clade-c** descendant from **clade-a** had a different set of co-mutation with Nsp3-ORF8 proteins.  
335 **Clade-d** descended further from **clade-c** had two mutation in Nsp13 (5768/5731) in addition  
336 to Nsp3/ORF8 proteins. Two strains from USA in the cluster radiating from **clade-d** harbored  
337 additional Nsp6 mutation stating them more divergent with scope of further possible  
338 evolution. The next subclade-e1 was found holding another new set of co-mutation of  
339 Nsp3/S/Nsp12. Whereas the highest number of co-mutations were found in subclade-e2 with  
340 combination of Nsp3/Nsp2/Nsp12/S/ORF3a prevalent in 30 genomes belonging to USA  
341 predicting them as active carrier of evolutionary force for SARS-CoV-2 divergence (Figure 3  
342 A & B). Presence of Nsp3 mutation (S1920P) in 238 strains underlined the origin of mutation  
343 from reference strain highlighting the first divergence in SARS-CoV-2 strain. In future, more  
344 and more genome availability from USA may indicate the evolutionary relationships with  
345 these co-mutations. Our result suggested that co-mutations are the major evolutionary force  
346 that drives the pathogenicity among the different geographical isolated strains which can  
347 responsible for higher and lower order of virulence among them.

#### 348 **The assessment of mutations in SARS-CoV-2 proteins**

349 Amino acid variations were predicted in eight (Nsp2, Nsp3, Nsp6, Nsp12, Nsp13, S, Orf3a,  
350 Orf8) SARS-CoV-2 proteins (Table 1). To identify their potential functional role, we carried  
351 out the structural analysis of the proteins. Pairwise sequence alignment of wild-type and  
352 mutant proteins provided the exact location and changes in amino acids. The GMQE and  
353 QMEAN values range from 0.45 to 0.72 and -1.43 to -2.81, respectively. The sequence  
354 identity ranges from 34% to 99%, which suggests that the models were constructed with high  
355 confidence and best quality (Figure 6). The I-Mutant DDG tool predicts if a mutation can  
356 largely destabilize the protein ( $\Delta\Delta G < -0.5$  Kcal/mol), largely stabilize ( $\Delta\Delta G > 0.5$  Kcal/mol) or  
357 have a weak effect ( $-0.5 \leq \Delta\Delta G \leq 0.5$  Kcal/mol). The protein stability analysis showed that all  
358 the identified mutations decreased the stability of seven proteins (Nsp2, Nsp6, Nsp12, Nsp13,  
359 S, Orf3a, Orf8) except Nsp3 (T1103P) which predicted to increase protein stability (Figure 6  
360 A-H). Further, to explore the role of mutations in SARS-CoV-2 proteins, we carried out  
361 HOPE analysis. D614G mutation in S-protein could disturb the rigidity of the protein and due  
362 to glycine, hydrophobicity will affect the intra hydrogen bond formation with G594. In ORF8  
363 and Nsp3, the mutation location was not conserved, so it did not affect or damage the protein

364 function. The mutation (P409L) in Nsp13 was present in the RNA virus helicase C-terminal  
365 domain. Since proline is a very rigid amino acid and therefore induce a particular backbone  
366 conformation that might be required at this position so this mutation could disturb domain and  
367 abolished its function. Mutation L37F (Nsp6) and T85I (Nsp2) were also highly conserved  
368 thus could profoundly damage the function of the respective protein. The P227L (Nsp12)  
369 mutation was in the RNA binding domain located on the surface of the protein; modification  
370 of this residue could disturb interactions with other molecules or other parts of the protein.  
371 Conclusively, Nsp3 mutation which appeared in all co-mutation combinations, contributed in  
372 increased protein stability among 238 strains could be assigned to their increased  
373 pathogenicity. Thus, we attempted to highlight the effects of these mutations in host pathogen  
374 interactions.

375

### 376 **Modelling of Host-Pathogen Interaction Network and its Functional Analysis**

377 The HPI Network of SARS-CoV-2 (HPIN-SARS-CoV-2) contained 58 edges, 56 nodes,  
378 including 5 viral and 51 host proteins (Figure 5). Number of degree (the number of edges per  
379 node) calculated based on HPI. The significant existence of few main gene hubs, namely, N, S  
380 and M in the network and the attraction of a large number of low-degree nodes toward each  
381 hub show strong evidence of controlling the topological properties of the network by these  
382 few hubs. N has 37 degrees, S, and M has 16 and 8 degrees, respectively. These viral proteins  
383 are the main hubs in the network, which regulate the network. Based on degree distribution,  
384 the viral protein N showed highest pathogenicity followed by S and M. N is a highly  
385 conserved major structural component of SARS-CoV virion involved in pathogenesis and  
386 used as a marker for diagnostic assays (Xia et al., 2020). Another structural protein S (spike  
387 glycoprotein), attach the virion to the cell membrane by interacting with host receptor,  
388 initiating the infection (Belouzard et al., 2012). The M protein, component of the viral  
389 envelope played a central role in virus morphogenesis and assembly via its interactions with  
390 other viral proteins (Garoff et al.,1998). Interestingly, we found four host proteins MYO5A,  
391 MYO5B, MYO5C and T had a maximum interaction with viral hub proteins. MYO5A,  
392 MYO5B, MYO5C interacting with all three (N, S and M) whereas T with two (S and M) viral  
393 hub proteins, showed a significant relationship with persistent infections caused by the SARS-  
394 CoV-2.

395 MYO5A, MYO5B and MYO5C proteins are Class V myosin (myosin-5) molecular motor that  
396 functions as an organelle transporter (Roland et al., 2011) (Sasaki, et al., 1995). It was found  
397 that the presence of myosin protein played a crucial role in coronavirus assembly and budding

398 in the infected cells (Neuman et al., 2008). These cytoskeletal proteins are of importance  
399 during internalization and subsequent intracellular transport of viral proteins. As we know at  
400 the entry level of virus, S interacts with host ACE2 receptor that internalizes the virus into the  
401 endosomes of the host cell inducing conformational changes in the S glycoprotein (Belouzard  
402 et al., 2012). It was found that inhibition of MYO5A, MYO5B, and MYO5C was efficient in  
403 blocking the internalization pathway, thus this target can be used for the development of a  
404 new treatment for SARS-CoV-2 (Dewerchin et al., 2014). Patients suffering from COVID-19  
405 undergo two major condition in the severe stage, thrombotic phenomenon and hypoxia, that  
406 are acting as silent killers (Bikdeli et al., 2020; Negri et al., 2020  
407 <https://doi.org/10.1101/2020.04.15.20067017>). Hypoxia, condition where oxygen level of the  
408 body reduces drastically results in the elevated expression of T protein in the body (Shao et  
409 al., 2015; Yoon et al., 2006). T protein (Brachyury/TBXT) is transcription factor involved in  
410 regulating genes required for mesoderm formation and differentiation thus playing an  
411 important role in pathogenesis. The detailed functional analysis of HPIN-SARS-CoV-2 was  
412 mapped on the radiological findings from the COVID-19 severely infected patients and non-  
413 survivors. It was reported that the levels of fibrin-degrading proteins, fibrinogen and D-dimer  
414 protein were 3-4 folds higher as compared to healthy individual. Therefore, reflecting  
415 coagulation activation from infection/sepsis, cytokine storm and impending multiple organs  
416 failure (Tang et al., 2020; Shi et al., 2020; Han et al., 2020, Li et al., 2020). In our network,  
417 we found 24 proteins (ANGPTL1, TNN, FGL2, ANGPTL6, TNC, FCN3, FCN2, ANGPTL4,  
418 FGB, FGA, ANGPT2, ANGPTL5, FGG, TNF, ANGPTL3, FCN1, FIBCD1, ANGPTL2,  
419 ANGPTL7, ANGPT4, MFAP4, FGL1, TNXB and ANGPT1) are associated with the above  
420 etiology (Figure 5 C). We also found the interaction of SMAD family proteins and SUMO1  
421 with N protein, which may lead to inhibition of apoptosis of infected lung cells (Zhao et al,  
422 2008). The interactome study reveals a significant role of identified host proteins in viral  
423 budding and related symptoms of COVID-19.

#### 424 **The mutation in SARS-CoV-2 proteins inhibit viral penetration into host**

425 In order to validate the effect of phenotypic variation (AAV), significant host proteins  
426 interactions from HPIN-SARS-CoV-2 were considered for *in silico* docking studies. Docking  
427 of S-Protein (wild type and mutant) with ACE2, TMPRSS2 and one of myosin proteins  
428 (MYO5C) were analyzed. Recent studies have shown that SARS-CoV-2 uses the ACE2 for  
429 entry and the serine protease TMPRSS2 for S protein priming (Wrapp et al., 2020). The poly-  
430 proteins (Nsp12, Nsp13, Nsp2, Nsp3 and Nsp6) of ORF1A and ORF1AB were docked with

431 RPS6 and ATP6V1G1 host proteins. The docking results showed that mutant S-protein could  
432 not bind efficiently with ACE2 and MYO5C, whereas mutation slightly promotes the binding  
433 with TMPRSS2 (Table 3, Figure 6 and Figure 5C). TMPRSS2 have been detected in both  
434 nasal and bronchial epithelium by immunohistochemistry (Bertram et. al., 2012), reported to  
435 occur largely in alveolar epithelial type II cells which are central to SARS-CoV-2  
436 pathogenesis (Furong et al., 2020). The wild-type S-protein form 16 hydrogen bonds and  
437 1058 non-bonded contacts with ACE2; whereas the mutant protein forms 12 hydrogen bond  
438 and 738 non-bonded contacts (Figure 6). This result suggests that D614G mutation in S-  
439 protein could affect viral entry into the host. Similarly, mutations present in the Nsp12,  
440 Nsp13, Nsp2, Nsp3 and Nsp6 of SARS-CoV-2 could inhibit the interaction with RPS6, but  
441 these mutations promote the binding with ATP6V1G1 expect Nsp6 (L3605F). The RPS6  
442 contributes to control cell growth and proliferation (Chauvin et al., 2014), so a loss of  
443 interaction with RPS6 could probably inhibit the production of viruses. Overall structural and  
444 interactome analyses suggests that identified mutations (Nsp2 (T265I), Nsp3 (S1920P), Nsp6  
445 (L3605F), Nsp12 (P4618L), Nsp13 (P5731L, Y5768C), S (D7611G)) in SARS-CoV-2 might  
446 play an important role in modifying the efficacy of viral entry and its pathogenesis. However,  
447 these observations required critical reevaluation as well as experimental work to confirm the  
448 *in-silico* results.

#### 449 **Regulation of SARS-CoV-2 pathogenicity by CpG island**

450 The genotyping analysis that we performed showed high frequency rate (45) of SNP at 5'UTR  
451 region (Table 1) and recent study also suggested that suppression of GC content could play a  
452 vital role in specific antiviral activities (Xia, 2020). As seen in SNP analysis, the common  
453 transitions of C>T and G>A alter the GC content of the SARS-CoV-2 (Table 1). This directed  
454 the analysis towards understanding the role of CpG island which is involved in silencing of  
455 transcription and down regulation of viral replication (Vivekanandan et al., 2010). Viral  
456 infections upregulate host DNA methyltransferase genes (DNMTs), and their overexpression  
457 leads to methylation of host CpG islands along with the viral CpGs (Vivekanandan et al.,  
458 2010). Since increased frequency of CpG motifs can serve as Pathogen-associated molecular  
459 pattern (PAMP) or Damage-associated molecular pattern (DAMP) which are potent inducers  
460 of strong innate immune responses (Barber, 2011; Frieman, 2008). Thus, CpG island profiling  
461 and their importance of existence in SARS-CoV-2 genomes was proceeded. We found that  
462 CpG islands were consistently present in two regions of the genome at the positions 285-385  
463 nucleotides (101 bp) and 28,324-28,425 nucleotides (102 bp). The results were consistent in

464 all 245 genomes analyzed in the present study with 100% conservancy in 237 genome  
465 sequences (Figure 7 A).

466 In the remaining 8 genomes, five genomes (MT246474.1 (G to A substitution at 354<sup>th</sup> position  
467 with respect to reference genome); MT276329.1, MT276330.1 and MT276598.1 (C to T  
468 substitution at 313<sup>th</sup> position) and MT246455.1 (G to T substitution at 332<sup>nd</sup> position)) showed  
469 point mutation in 5' CpG island; whereas three genomes (MT159718.1 (C to T substitution at  
470 28409<sup>th</sup> position); MT159717.1 and MT184911.1 (G to T substitution at 28378<sup>th</sup> position))  
471 showed point mutation in 3'CpG end (Figure 7 D). Interestingly, all these sequences belong to  
472 USA. On further locating CpG island positions with respect to proteins, it was found that  
473 these two CpG islands were located at two prime locations within the genome, one in Nsp1  
474 (Figure 7 B), and another within Nucleocapsid (N) protein (Figure 7 C). Previously, it was  
475 reported that both the proteins interacted with 5' UTR region playing crucial roles in viral  
476 replication and gene expressions (Guan et al., 2012; Yang and Leibowitz, 2015; Galan et al.,  
477 2005). Most pivotal role of N protein revolves around encapsulation of viral gRNA which  
478 leads to formation of ribonucleoprotein complex (RNP), which is a vital step in assembly of  
479 viral particles (Cong et al., 2017).

480 Nsp1 protein in coronaviruses plays a regulatory role in transcription and viral replication  
481 (Cong et al., 2017). It is known to interact with 5' UTR of host cell mRNA to induce its  
482 endonucleolytic cleavage (Huang et al., 2011; Narayanan et al., 2015), thus inhibiting host  
483 gene expression (Kamitani et al., 2009). It also plays an important role in blocking IFN-  
484 dependent antiviral signaling pathways leading to dysregulation of host immune system  
485 (Kamitani et al., 2006; Wathelet et al., 2007; Law et al., 2007). CpG sites can be targeted by  
486 Zinc Finger Antiviral Proteins which can mediate antiviral restriction through CpG motif  
487 detection (Bick et al, 2003; Liu et al., 2015; Chiu et al., 2018). Apart from this, CpG  
488 oligodeoxynucleotides (ODNs) are known to act as adjuvants and are already established as a  
489 potent stimulator for host immune system (Campbell, 2017; Becker, 2005; Yuan, 2017; Singh  
490 et al., 2016; Yu et al., 2018). Moreover, recent studies conducted on influenza A genome and  
491 Zika virus genome has shown that by increasing the CpG dinucleotides in viral genome,  
492 impairment of viral infection is observed (Gaunt et al., 2016; Trus et al., 2020). Our result  
493 showed that the presence of conserved CpG islands in Nps1 and N protein across 245  
494 genomes of SARS-CoV-2 indicated their role in pathogenesis and can be targeted by Zinc  
495 Finger Antiviral Proteins or exploited to design CpG-recoded vaccines.

## 496 **Conclusions**

497 The genomic and proteomic survey of 245 SARS-CoV-2 strains reported from subset of  
498 population of different countries reflected global transmission during the outbreak of COVID-  
499 19. The viral phylogenetic network with five clads (a-e) provided a landscape of the current  
500 stage of epidemic where major divergence was observed in USA strains. From this we  
501 propose genotypes linked to geographic clads in which signature SNP can be used to track and  
502 monitor the epidemic. Demarcation of co-mutation in the SARS-CoV-2 strains by assessing  
503 co-mutations also highlighted the evolutionary relationships among the viral proteins. Our  
504 results suggested that co-mutation are indicative of AAV based induced pathogenicity leading  
505 to multiple mutations embedded in few genomes. However, co-mutations are still in  
506 evolutionary process and more combinations can be predicted with a large dataset. High-  
507 frequency AAV mutations were present in the critical proteins, including the Nsp2, Nsp3,  
508 Nsp6, Nsp12, Nsp13, S, Orf3a, Orf8 which could be considered for designing a vaccine.  
509 Comparative analysis of proteins from wild and mutated strains showed positive selection of  
510 mutation in Nsp3 but not in rest of the mutants. HPI model can be used as the fundamental  
511 basis for structure-guided pathogenesis process inside host cell. The interactome study  
512 showed MYO-5 proteins as a key host partner and also highlighted the key role of N, S and M  
513 viral proteins for conferring SARS-CoV-2 pathogenicity. The mutation in the S protein could  
514 affect the viral entry by loose binding with ACE2. The presence of CpG islands in N and  
515 Nsp1 protein could play a critical role in pathogenesis regulation. Based on our multi-omics  
516 approach: genomics, proteomics, interactomics and structural biology; provided an  
517 opportunity for better understanding of COVID-19 pandemic and can be considered in  
518 ongoing vaccine development programs.

## 519 **Authors Contribution**

520 RL, VG, SH, MV conceived and designed the study. VG, HV, SH, NS, KP, MZM performed  
521 the analysis and develop figures. VG, SH, MV, KP wrote the manuscript and RL, RK, HV,  
522 US, PH, SS help in shaping manuscript.

## 523 **Conflict of Interest**

524 Authors declare no conflict of interest

## 525 **Acknowledgements**



526 VG acknowledge Phixgen Pvt. Ltd. for research fellowship. MV, SS acknowledge Dr. P.  
527 Hemalatha Reddy, Principal, Sri Venkateswara College, University of Delhi for her constant  
528 support and encouragement. RL and US also acknowledge The National Academy of  
529 Sciences, India, for support under the NASI-Senior Scientist Platinum Jubilee Fellowship  
530 Scheme. NS acknowledge Council of Scientific and Industrial Research (CSIR), New Delhi  
531 for doctoral fellowships. KP thanks Hub of Bioinformatics for providing support. SH would  
532 like to thank Jaypee Institute of Information Technology, Noida India for providing support.  
533 HV would like to thank Ramjas College, University of Delhi, Delhi for providing support. RK  
534 acknowledges Magadh University, Bodh Gaya for providing support. MZM acknowledge  
535 Department of Health Welfare, Government of India under young scientist scheme for  
536 financial support. PH would like to thank Maitreyi College, University of Delhi, Delhi for  
537 providing support.

538

539

540

541

## 542 **References**

- 543 • Abduljalil, J. M. & Abduljalil, B. M. Epidemiology, genome, and clinical features of  
544 the pandemic SARS-CoV-2: a recent view. *New Microbes New Infect* **35**, 100672-  
545 100672, doi:10.1016/j.nmni.2020.100672 (2020).
- 546 • Ammari, M. G., Gresham, C. R., McCarthy, F. M. & Nanduri, B. HPIDB 2.0: a  
547 curated database for host–pathogen interactions. *Database* **2016**,  
548 doi:10.1093/database/baw103 (2016).
- 549 • Assenov, Y., Ramírez, F., Schelhorn, S.-E., Lengauer, T. & Albrecht, M. Computing  
550 topological parameters of biological networks. *Bioinformatics* **24**, 282-284,  
551 doi:10.1093/bioinformatics/btm554 %J Bioinformatics (2007).
- 552 • Astuti, I. & Ysrafil. Severe Acute Respiratory Syndrome Coronavirus 2 (SARS-CoV-  
553 2): An overview of viral structure and host response. *Diabetes Metab Syndr* **14**, 407-  
554 412, doi:10.1016/j.dsx.2020.04.020 (2020).
- 555 • Barber, G. N. Cytoplasmic DNA innate immune pathways. **243**, 99-108,  
556 doi:10.1111/j.1600-065X.2011.01051.x (2011).

- 557 • Becker, Y. CpG ODNs treatments of HIV-1 infected patients may cause the decline of  
558 transmission in high risk populations - a review, hypothesis and implications. *Virus*  
559 *Genes* **30**, 251-266, doi:10.1007/s11262-004-5632-2 (2005).
- 560 • Begum, F., Banerjee, A. K., Tripathi, P. P. & Ray, U. Two mutations P/L and Y/C in  
561 SARS-CoV-2 helicase domain exist together and influence helicase RNA binding.  
562 *bioRxiv*, 2020.2005.2014.095224, doi:10.1101/2020.05.14.095224 (2020).
- 563 • Begum, F. *et al.* Analyses of spike protein from first deposited sequences of SARS-  
564 CoV2 from West Bengal, India. *bioRxiv*, 2020.2004.2028.066985,  
565 doi:10.1101/2020.04.28.066985 (2020).
- 566 • Belouzard, S., Millet, J. K., Licitra, B. N. & Whittaker, G. R. Mechanisms of  
567 Coronavirus Cell Entry Mediated by the Viral Spike Protein. **4**, 1011-1033 (2012).
- 568 • Benvenuto, D. *et al.* Evolutionary analysis of SARS-CoV-2: how mutation of Non-  
569 Structural Protein 6 (NSP6) could affect viral autophagy. *Journal of Infection*,  
570 doi:https://doi.org/10.1016/j.jinf.2020.03.058 (2020).
- 571 • Bertram, S. *et al.* Influenza and SARS-coronavirus activating proteases TMPRSS2 and  
572 HAT are expressed at multiple sites in human respiratory and gastrointestinal tracts.  
573 *PLoS One* **7**, e35876-e35876, doi:10.1371/journal.pone.0035876 (2012).
- 574 • Bhattacharyya, C. *et al.* Global Spread of SARS-CoV-2 Subtype with Spike Protein  
575 Mutation D614G is Shaped by Human Genomic Variations that Regulate Expression  
576 of <em>TMPRSS2</em> and <em>MX1</em>. *Genes*.  
577 *bioRxiv*, 2020.2005.2004.075911, doi:10.1101/2020.05.04.075911 (2020).
- 578 • Bick, M. J. *et al.* Expression of the Zinc-Finger Antiviral Protein Inhibits Alphavirus  
579 Replication. *Journal of Virology* **77**, 11555, doi:10.1128/JVI.77.21.11555-11562.2003  
580 (2003).
- 581 • Bikdeli, B. *et al.* COVID-19 and Thrombotic or Thromboembolic Disease:  
582 Implications for Prevention, Antithrombotic Therapy, and Follow-up. *Journal of the*  
583 *American College of Cardiology*, 27284, doi:10.1016/j.jacc.2020.04.031 (2020).
- 584 • Campbell, J. D. Development of the CpG Adjuvant 1018: A Case Study. *Methods Mol*  
585 *Biol* **1494**, 15-27, doi:10.1007/978-1-4939-6445-1\_2 (2017).
- 586 • Capriotti, E., Calabrese, R. & Casadio, R. Predicting the insurgence of human genetic  
587 diseases associated to single point protein mutations with support vector machines and  
588 evolutionary information. *Bioinformatics* **22**, 2729-2734,  
589 doi:10.1093/bioinformatics/btl423 %J Bioinformatics (2006).

- 590 • Chauvin, C. *et al.* Ribosomal protein S6 kinase activity controls the ribosome  
591 biogenesis transcriptional program. *Oncogene* **33**, 474-483, doi:10.1038/onc.2012.606  
592 (2014).
- 593 • Cheng, S.-C. *et al.* First case of Coronavirus Disease 2019 (COVID-19) pneumonia in  
594 Taiwan. *Journal of the Formosan Medical Association* **119**, 747-751,  
595 doi:<https://doi.org/10.1016/j.jfma.2020.02.007> (2020).
- 596 • Chiu, H.-P. *et al.* Inhibition of Japanese encephalitis virus infection by the host zinc-  
597 finger antiviral protein. *PLOS Pathogens* **14**, e1007166,  
598 doi:10.1371/journal.ppat.1007166 (2018).
- 599 • Cong, Y., Kriegenburg, F., de Haan, C. A. M. & Reggiori, F. Coronavirus  
600 nucleocapsid proteins assemble constitutively in high molecular oligomers. *Scientific*  
601 *Reports* **7**, 5740, doi:10.1038/s41598-017-06062-w (2017).
- 602 • Cong, Y. *et al.* Nucleocapsid Protein Recruitment to Replication-Transcription  
603 Complexes Plays a Crucial Role in Coronaviral Life Cycle. *Journal of virology* **94**,  
604 e01925-01919, doi:10.1128/JVI.01925-19 (2020).
- 605 • Cottam, E. M., Whelband, M. C. & Wileman, T. Coronavirus NSP6 restricts  
606 autophagosome expansion. *Autophagy* **10**, 1426-1441, doi:10.4161/auto.29309 (2014).
- 607 • Dewerchin, H. L., Desmarests, L. M., Noppe, Y. & Nauwynck, H. J. Myosins 1 and 6,  
608 myosin light chain kinase, actin and microtubules cooperate during antibody-mediated  
609 internalisation and trafficking of membrane-expressed viral antigens in feline  
610 infectious peritonitis virus infected monocytes. *Vet Res* **45**, 17-17, doi:10.1186/1297-  
611 9716-45-17 (2014).
- 612 • Forster, P., Forster, L., Renfrew, C. & Forster, M. Phylogenetic network analysis of  
613 SARS-CoV-2 genomes. *Proceedings of the National Academy of Sciences* **117**, 9241,  
614 doi:10.1073/pnas.2004999117 (2020).
- 615 • Frieman, M., Heise, M. & Baric, R. SARS coronavirus and innate immunity. *Virus*  
616 *Research* **133**, 101-112, doi:<https://doi.org/10.1016/j.virusres.2007.03.015> (2008).
- 617 • Galán, C., Enjuanes, L. & Almazán, F. A point mutation within the replicase gene  
618 differentially affects coronavirus genome versus minigenome replication. *Journal of*  
619 *virology* **79**, 15016-15026, doi:10.1128/JVI.79.24.15016-15026.2005 (2005).
- 620 • Garoff, H., Hewson, R. & Opstelten, D. J. Virus maturation by budding. *Microbiol*  
621 *Mol Biol Rev* **62**, 1171-1190 (1998).

- 622 • Gaunt, E. *et al.* Elevation of CpG frequencies in influenza A genome attenuates  
623 pathogenicity but enhances host response to infection. *Elife* **5**, e12735-e12735,  
624 doi:10.7554/eLife.12735 (2016).
- 625 • Guan, B.-J., Su, Y.-P., Wu, H.-Y. & Brian, D. A. Genetic evidence of a long-range  
626 RNA-RNA interaction between the genomic 5' untranslated region and the  
627 nonstructural protein 1 coding region in murine and bovine coronaviruses. *Journal of*  
628 *virology* **86**, 4631-4643, doi:10.1128/JVI.06265-11 (2012).
- 629 • Han, H. *et al.* Prominent changes in blood coagulation of patients with SARS-CoV-2  
630 infection. *Clin Chem Lab Med*, doi:10.1515/cclm-2020-0188 (2020).
- 631 • Hassan, S. S. *et al.* On spatial molecular arrangements of SARS-CoV2 genomes of  
632 Indian patients. *bioRxiv*, 2020.2005.2001.071985, doi:10.1101/2020.05.01.071985  
633 (2020).
- 634 • Hoelzer, K., Shackelton, L. A. & Parrish, C. R. Presence and role of cytosine  
635 methylation in DNA viruses of animals. *Nucleic Acids Research* **36**, 2825-2837,  
636 doi:10.1093/nar/gkn121 %J Nucleic Acids Research (2008).
- 637 • Huang, C. *et al.* SARS Coronavirus nsp1 Protein Induces Template-Dependent  
638 Endonucleolytic Cleavage of mRNAs: Viral mRNAs Are Resistant to nsp1-Induced  
639 RNA Cleavage. *PLOS Pathogens* **7**, e1002433, doi:10.1371/journal.ppat.1002433  
640 (2011).
- 641 • Jacobson, M. P., Friesner, R. A., Xiang, Z. & Honig, B. On the Role of the Crystal  
642 Environment in Determining Protein Side-chain Conformations. *Journal of Molecular*  
643 *Biology* **320**, 597-608, doi:https://doi.org/10.1016/S0022-2836(02)00470-9 (2002).
- 644 • Jones, P. A. Functions of DNA methylation: islands, start sites, gene bodies and  
645 beyond. *Nature Reviews Genetics* **13**, 484-492, doi:10.1038/nrg3230 (2012).
- 646 • Kamitani, W., Huang, C., Narayanan, K., Lokugamage, K. G. & Makino, S. A two-  
647 pronged strategy to suppress host protein synthesis by SARS coronavirus Nsp1  
648 protein. *Nat Struct Mol Biol* **16**, 1134-1140, doi:10.1038/nsmb.1680 (2009).
- 649 • Kamitani, W. *et al.* Severe acute respiratory syndrome coronavirus nsp1 protein  
650 suppresses host gene expression by promoting host mRNA degradation. *Proc Natl*  
651 *Acad Sci U S A* **103**, 12885-12890, doi:10.1073/pnas.0603144103 (2006).
- 652 • Kelley, L. A., Mezulis, S., Yates, C. M., Wass, M. N. & Sternberg, M. J. E. The  
653 Phyre2 web portal for protein modeling, prediction and analysis. *Nature Protocols* **10**,  
654 845-858, doi:10.1038/nprot.2015.053 (2015).

- 655 • Khailany, R. A., Safdar, M. & Ozaslan, M. Genomic characterization of a novel  
656 SARS-CoV-2. *Gene Reports* **19**, 100682,  
657 doi:<https://doi.org/10.1016/j.genrep.2020.100682> (2020).
- 658 • Koyama, T., Platt, D. & Parida, L. *Variant analysis of COVID-19 genomes.* (2020).
- 659 • Krinner, S. *et al.* CpG domains downstream of TSSs promote high levels of gene  
660 expression. *Nucleic Acids Research* **42**, 3551-3564, doi:10.1093/nar/gkt1358 %J  
661 Nucleic Acids Research (2014).
- 662 • Law, A. H. Y., Lee, D. C. W., Cheung, B. K. W., Yim, H. C. H. & Lau, A. S. Y. Role  
663 for Nonstructural Protein 1 of Severe Acute Respiratory Syndrome Coronavirus in  
664 Chemokine Dysregulation. *Journal of virology* **81**, 2537-2537, doi:10.1128/jvi.02744-  
665 06 (2007).
- 666 • Lei, J., Kusov, Y. & Hilgenfeld, R. Nsp3 of coronaviruses: Structures and functions of  
667 a large multi-domain protein. *Antiviral Res* **149**, 58-74,  
668 doi:10.1016/j.antiviral.2017.11.001 (2018).
- 669 • Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL): an online tool for phylogenetic  
670 tree display and annotation. *Bioinformatics* **23**, 127-128,  
671 doi:10.1093/bioinformatics/btl529 %J Bioinformatics (2006).
- 672 • Li, Q. *et al.* Early Transmission Dynamics in Wuhan, China, of Novel Coronavirus-  
673 Infected Pneumonia. *The New England journal of medicine* **382**, 1199-1207,  
674 doi:10.1056/NEJMoa2001316 (2020).
- 675 • Li, T., Lu, H. & Zhang, W. Clinical observation and management of COVID-19  
676 patients. *Emerging Microbes & Infections* **9**, 687-690,  
677 doi:10.1080/22221751.2020.1741327 (2020).
- 678 • Liu, C.-H., Zhou, L., Chen, G. & Krug, R. M. Battle between influenza A virus and a  
679 newly identified antiviral activity of the PARP-containing ZAPL protein. *Proceedings*  
680 *of the National Academy of Sciences*, 201509745, doi:10.1073/pnas.1509745112  
681 (2015).
- 682 • Mandal, S., Singh, R.S., Sharma, S.K., Malik, M.Z. and Singh, R.B. Complexity in  
683 SARS-CoV-2 genome data: Price theory of mutant isolates. *bioRxiv*,  
684 doi.org/10.1101/2020.05.04.077511 (2020).
- 685 • Nafis S, Kalaiarasan P, Brojen Singh RK, Husain M, Bamezai RN. Apoptosis  
686 regulatory protein-protein interaction demonstrates hierarchical scale-free fractal  
687 network. *Brief Bioinform.* 2015;16(4):675-699. doi:10.1093/bib/bbu036

- 688 • Narayanan, K., Ramirez, S. I., Lokugamage, K. G. & Makino, S. Coronavirus  
689 nonstructural protein 1: Common and distinct functions in the regulation of host and  
690 viral gene expression. *Virus Res* **202**, 89-100, doi:10.1016/j.virusres.2014.11.019  
691 (2015).
- 692 • Neuman, B. W. *et al.* Proteomics analysis unravels the functional repertoire of  
693 coronavirus nonstructural protein 3. *Journal of virology* **82**, 5279-5294,  
694 doi:10.1128/JVI.02631-07 (2008).
- 695 • NHS Press Conference, Feb. 4 2020 - National Health Commission (NHC) of the  
696 People's Republic of China
- 697 • Pachetti, M. *et al.* Emerging SARS-CoV-2 mutation hot spots include a novel RNA-  
698 dependent-RNA polymerase variant. *Journal of Translational Medicine* **18**,  
699 doi:10.1186/s12967-020-02344-6 (2020).
- 700 • Pyrc, K., Berkhout, B. & van der Hoek, L. The Novel Human Coronaviruses NL63  
701 and HKU1. *Journal of Virology* **81**, 3051, doi:10.1128/JVI.01466-06 (2007).
- 702 • Qi, F., Qian, S., Zhang, S. & Zhang, Z. Single cell RNA sequencing of 13 human  
703 tissues identify cell types and receptors of human coronaviruses. *Biochemical and*  
704 *Biophysical Research Communications* **526**, 135-140,  
705 doi:<https://doi.org/10.1016/j.bbrc.2020.03.044> (2020).
- 706 • Ramachandran, S., Kota, P., Ding, F. & Dokholyan, N. V. Automated minimization of  
707 steric clashes in protein structures. *Proteins* **79**, 261-270, doi:10.1002/prot.22879  
708 (2011).
- 709 • Roland, J. T. *et al.* Rab GTPase–Myo5B complexes control membrane recycling and  
710 epithelial polarization. *Proceedings of the National Academy of Sciences* **108**, 2789,  
711 doi:10.1073/pnas.1010754108 (2011).
- 712 • Sasaki, H. *et al.* Myosin-actin interaction plays an important role in human  
713 immunodeficiency virus type 1 release from host cells. *Proc Natl Acad Sci U S A* **92**,  
714 2026-2030, doi:10.1073/pnas.92.6.2026 (1995).
- 715 • Schneidman-Duhovny, D., Inbar, Y., Nussinov, R. & Wolfson, H. J. PatchDock and  
716 SymmDock: servers for rigid and symmetric docking. *Nucleic acids research* **33**,  
717 W363-W367, doi:10.1093/nar/gki481 (2005).
- 718 • Sexton, N. R. *et al.* Homology-Based Identification of a Mutation in the Coronavirus  
719 RNA-Dependent RNA Polymerase That Confers Resistance to Multiple Mutagens.  
720 *Journal of virology* **90**, 7415-7428, doi:10.1128/JVI.00080-16 (2016).

- 721 • Shi, H. *et al.* Radiological findings from 81 patients with COVID-19 pneumonia in  
722 Wuhan, China: a descriptive study. *Lancet Infect Dis* **20**, 425-434, doi:10.1016/S1473-  
723 3099(20)30086-4 (2020).
- 724 • Shiraishi, M., Sekiguchi, A., Oates, A. J., Terry, M. J. & Miyamoto, Y. HOX gene  
725 clusters are hotspots of de novo methylation in CpG islands of human lung  
726 adenocarcinomas. *Oncogene* **21**, 3659-3662, doi:10.1038/sj.onc.1205453 (2002).
- 727 • Sievers, F. *et al.* Fast, scalable generation of high-quality protein multiple sequence  
728 alignments using Clustal Omega. *Mol Syst Biol* **7**, 539-539, doi:10.1038/msb.2011.75  
729 (2011).
- 730 • Singh, S. M. *et al.* Characterization of Immune Responses to an Inactivated Avian  
731 Influenza Virus Vaccine Adjuvanted with Nanoparticles Containing CpG ODN. *Viral*  
732 *Immunology* **29**, 269-275, doi:10.1089/vim.2015.0144 (2016).
- 733 • Sharma, S., Singh, I., Haider, S., Malik, M.Z., Ponnusamy, K. and Rai, E.. ACE  
734 Homo-dimerization, Human Genomic variants and Interaction of Host Proteins  
735 Explain High Population Specific Differences in Outcomes of COVID19. bioRxiv.  
736 doi.org/10.1101/2020.04.24.050534 (2020).
- 737 • Tang, X. *et al.* On the origin and continuing evolution of SARS-CoV-2. *National*  
738 *Science Review*, doi:10.1093/nsr/nwaa036 (2020).
- 739 • Trus, I. *et al.* CpG-Recoding in Zika Virus Genome Causes Host-Age-Dependent  
740 Attenuation of Infection With Protection Against Lethal Heterologous Challenge in  
741 Mice. **10**, doi:10.3389/fimmu.2019.03077 (2020).
- 742 • Venselaar, H., Te Beek, T. A. H., Kuipers, R. K. P., Hekkelman, M. L. & Vriend, G.  
743 Protein structure analysis of mutations causing inheritable diseases. An e-Science  
744 approach with life scientist friendly interfaces. *BMC Bioinformatics* **11**, 548-548,  
745 doi:10.1186/1471-2105-11-548 (2010).
- 746 • Vivekanandan, P., Daniel, H. D., Kannangai, R., Martinez-Murillo, F. & Torbenson,  
747 M. Hepatitis B virus replication induces methylation of both host and viral DNA. *J*  
748 *Virol* **84**, 4321-4329, doi:10.1128/jvi.02280-09 (2010).
- 749 • Waterhouse, A. *et al.* SWISS-MODEL: homology modelling of protein structures and  
750 complexes. *Nucleic acids research* **46**, W296-W303, doi:10.1093/nar/gky427 (2018).
- 751 • Wathelet, M. G., Orr, M., Frieman, M. B. & Baric, R. S. Severe acute respiratory  
752 syndrome coronavirus evades antiviral signaling: role of nsp1 and rational design of an

- 753 attenuated strain. *Journal of Virology* **81**, 11620-11633, doi:10.1128/jvi.00702-07  
754 (2007).
- 755 • Wrapp, D. *et al.* Structural Basis for Potent Neutralization of Betacoronaviruses by  
756 Single-Domain Camelid Antibodies. *Cell* **181**, 1004-1015.e1015,  
757 doi:<https://doi.org/10.1016/j.cell.2020.04.031> (2020).
  - 758 • Wu, C. *et al.* Analysis of therapeutic targets for SARS-CoV-2 and discovery of  
759 potential drugs by computational methods. *Acta Pharmaceutica Sinica B*,  
760 doi:<https://doi.org/10.1016/j.apsb.2020.02.008> (2020).
  - 761 • Xia, X. Extreme Genomic CpG Deficiency in SARS-CoV-2 and Evasion of Host  
762 Antiviral Defense. *Molecular Biology and Evolution*, doi:10.1093/molbev/msaa094  
763 (2020).
  - 764 • Yang, D. & Leibowitz, J. L. The structure and functions of coronavirus genomic 3' and  
765 5' ends. *Virus Research* **206**, 120-133,  
766 doi:<https://doi.org/10.1016/j.virusres.2015.02.025> (2015).
  - 767 • Yoon, D. *et al.* Hypoxia-inducible Factor-1 Deficiency Results in Dysregulated  
768 Erythropoiesis Signaling and Iron Homeostasis in Mouse Development. **281**, 25703-  
769 25711, doi:10.1074/jbc.M602329200 (2006).
  - 770 • Yu, C.-H., Qin, Z., Martin-Martinez, F. J. & Buehler, M. J. A Self-Consistent  
771 Sonification Method to Translate Amino Acid Sequences into Musical Compositions  
772 and Application in Protein Design Using Artificial Intelligence. *ACS Nano* **13**, 7471-  
773 7482, doi:10.1021/acsnano.9b02180 (2019).
  - 774 • Yu, P. *et al.* A CpG oligodeoxynucleotide enhances the immune response to rabies  
775 vaccination in mice. *Virol J* **15**, 174-174, doi:10.1186/s12985-018-1089-1 (2018).
  - 776 • Yuan, F. *et al.* Immunoprotection induced by CpG-ODN/Poly(I:C) combined with  
777 recombinant gp90 protein in chickens against reticuloendotheliosis virus infection.  
778 *Antiviral Res* **147**, 1-10, doi:<https://doi.org/10.1016/j.antiviral.2017.04.019> (2017)

779

## 780 Tables

781 **Table 1:** Common SNP and AAV mutations occurring in SARS CoV-2 genomes

782

CDS	Point Mutation	Position	Frequency	Amino Acid-Residue	Variant	Position	Frequency
-----	----------------	----------	-----------	--------------------	---------	----------	-----------



5'UTR	C→T	67	45				
Nsp2	C→T	885	29	T	I	265	31
Nsp2	C→T	2863	44				
Nsp3/PL-PRO	T→C	5852	238	S	P	1920	238
Nsp4	C→T	8608	88				
Nsp6	G→T	10909	21	L	F	3605	21
Nsp8	T→C	12299	238				
Nsp12 (RdRp)	C→T	14234	44	P	L	4618	46
Nsp13 (Hel)	C→T	17573	63	P	L	5731	64
Nsp13 (Hel)	A→G	17684	63	Y	C	5768	64
Nsp13 (Hel)	C→T	17886	68				
S	A→G	23232	43	D	G	7611	45
Orf3a	G→T	25392	32	Q	H	8327	34
Orf8	T→C	27973	88	L	S	9033	89

783

784 **Table 2:** Co-mutations combinations and genomic location identified in different proteins of  
 785 SARS-COV-2

Positions	Variation(s)	(Co)Mutations	Mutated protein	Descendants
1920	S>P	Nsp3	1	87
5768/5731/1920/9033	Y>C/P>L/S>P/L>S	Nsp13_1/Nsp13_2/Nsp3/ORF8	4	62
1920/9033	S>P/L>S	Nsp3-ORF8	1	22
4618/7611/8327/1920/265	P>L/D>G/Q>H/S>P/T>I	nsp12/S/ORF3a/Nsp3/Nsp2	5	30
4618/8327/1920/265	P>L/Q>H/S>P/T>I	Nsp12/ORF3a/Nsp3/Nsp2	4	1
4618/7611/8327/1920	P>L/D>G/Q>H/S>P	Nsp12/S/ORF3a/Nsp3	4	3
3605/1920	L>F/S>P	Nsp6/Nsp3	2	16
3605/1920/9033	L>F/S>P/L>S	Nsp6/Nsp3/ORF8	3	3
5768/5731/3605/1920/9033	Y>C/P>L/L>F/S>P/L>S	Nsp13_1/Nsp13_2/Nsp6/Nsp3/ORF8	5	2

4618/7611/1920	P>L/D>G/S>P	Nsp12/S/Nsp3	3	12
----------------	-------------	--------------	---	----

786

787 **Table 3.** *In silico* docking analysis of SARS-CoV-2 proteins with Human proteins

SARS CoV-2	Host Protein	Wild type score	Mutant score	Difference*
S-Protein	ACE2	18296	17722	574
S-Protein	TRMPSS2	20284	21180	-896
S-Protein	MYO5C	18538	17390	1148
Nsp13	RPS6	17772	15750	2022
Nsp13	ATP6V1G1	14432	20242	-5810
Nsp12	RPS6	16570	15750	820
Nsp12	ATP6V1G1	17150	20242	-3092
Nsp6	RPS6	19336	17736	1600
Nsp6	ATP6V1G1	17614	16022	1592
Nsp3	RPS6	22888	21866	1022
Nsp3	ATP6V1G1	20760	21070	-310
Nsp2	RPS6	22584	19540	3044
Nsp2	ATP6V1G1	18402	18592	-190

788

789

790

791

## 792 Figure Legends

793 **Figure 1.** Phylogenetic network of 245 SARS-CoV-2 genomes. (A) Nucleotide based  
 794 phylogenetic analysis of SARS-CoV-2 isolates using the Maximum Likelihood method based  
 795 on the Tamura-Nei model, (B) Amino acid based phylogenomic analysis. Circle areas are  
 796 proportional to the number of taxa. The map is diverged into 5 major clade (a-e) representing  
 797 variation in the genomes at amino-acid level. The colored circle represents the country of  
 798 origin of each isolate.

799 **Figure 2.** Distribution of SNP (A, B) and AAV (C, D) mutations of SARS-CoV-2 isolates  
 800 from the globe. (A) Frequency based plot of 12 possible SNP mutations across 245 genomes,  
 801 (B) Frequencies of the single SNP mutations with locations on the genome, (C) AAV based  
 802 mutations across the genomes, (D) Top 9 AAV mutations holding highest frequencies among  
 803 245 genomes and their respective positions. The nucleotide and amino-acid positions are  
 804 based on the reference genome of SARS-CoV-2.

805 **Figure 3.** (A) AAV based phylogenetic map of 245 SARS-CoV-2 genomes. Node color  
 806 represents co-mutational combinations. The formation of each clade is well correlated with  
 807 the mutational combinations (n=10). (B) Genetic separation among the SARS-CoV-2 strains  
 808 showing divergent evolution (from ancestral node a to b, c & e) and directive evolution (node  
 809 c to d) with adjoining daughter nodes represented by #.

810 **Figure 4.** 3-D structure prediction of SARS-CoV-2 proteins harboring mutations at different  
811 locations to check for its stability in the cell. Structure are predicted using SwissModel and  
812 Phyre2 servers.

813 **Figure 5.** (A) Host-pathogenic interaction map of SARS-CoV-2 and human proteins. Red  
814 triangles represent viral proteins found to be directly interacting with the human proteins,  
815 whereas the pink nodes denote the human proteins. Four major hubs were identified (green)  
816 found interacting with maximum viral proteins. (B) Number of degree (the number of edges  
817 per node) calculated based on HPI. The significant existence of few main gene hubs, namely,  
818 N, S and M in the network and the attraction of a large number of low-degree nodes toward  
819 each hub show strong evidence of controlling the topological properties of the network by  
820 these few hubs. N has 37 degrees, S, and M has 16 and 8 degrees, respectively. These viral  
821 proteins are the main hubs in the network, which regulate the network. (C) Functional  
822 enrichment of the human proteins (grey nodes) found to be directly interacting with viral  
823 proteins (pink nodes). Color bar around the nodes represent their functionality role in human  
824 body.

825

826 **Figure 6.** *In-silico* receptor-ligand docking analysis for mutated S protein (D7611G) from  
827 SARS-CoV-2 and ACE2 protein present in human. B & C represents amino-acid interactions  
828 between wild type and mutated Spike protein with ACE2 receptor.

829

830 **Figure 7.** (A) Detection of two CpG islands in Wuhan\_Hu-1 complete genome sequence  
831 (Accession number: MT121215.1), marked by blue arrows with their respective positions. (B)  
832 One of the CpG island was found to be located towards the 5' end of the genome, in ORF1ab.  
833 (C) Another CpG island was found towards the 3' end of the genome, located in ORF9 coding  
834 for N protein. (D) Five strains of USA showing point mutations in CpG island 1 located on 5'  
835 end of genome (at positions 313, 332 and 354 with respect to reference genome) and Three  
836 sequences showing substitutions in CpG island 2 at 28367, 28378 and 28409 positions  
837 respectively.

838

839

840

841

842

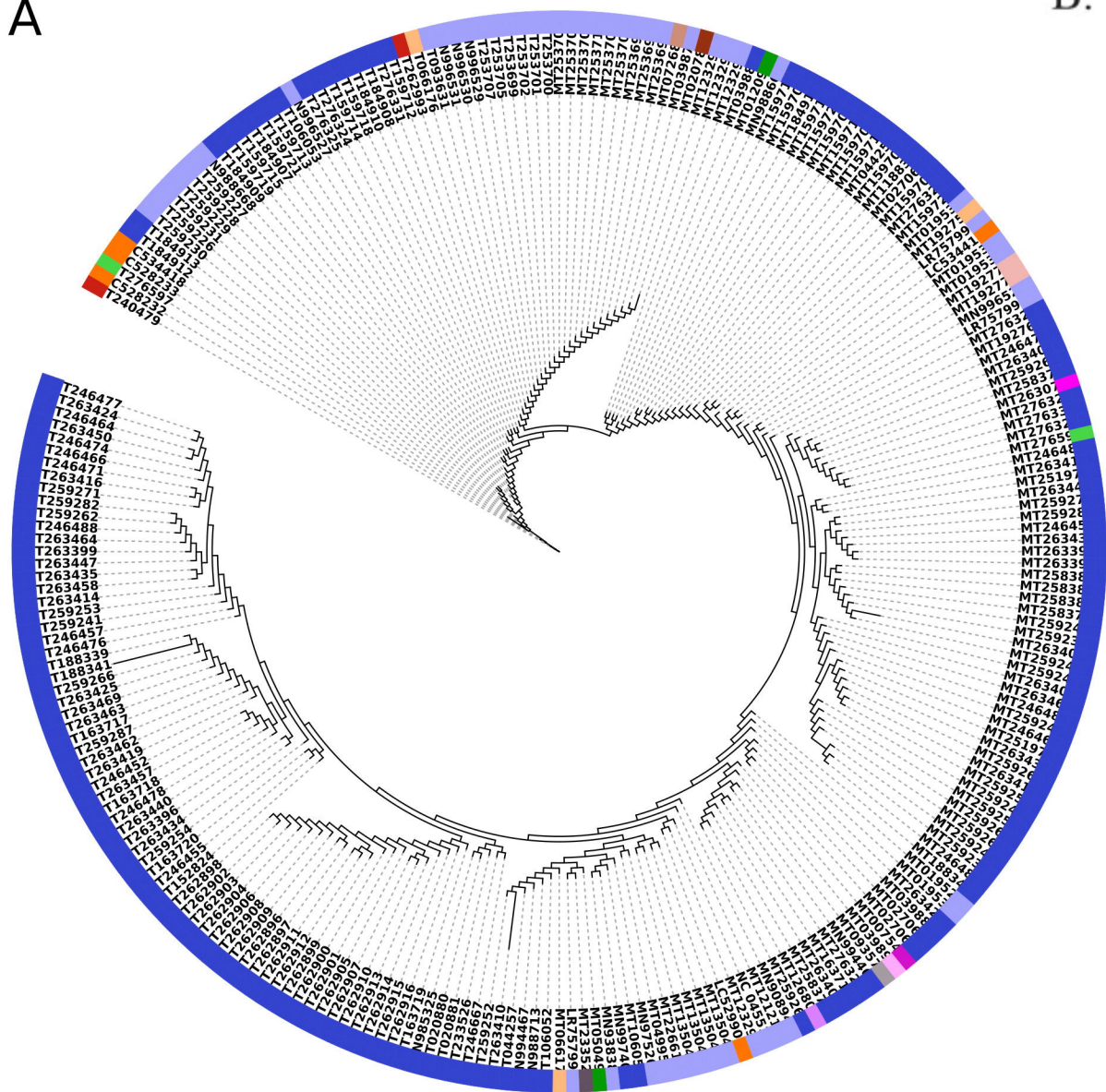
843

844

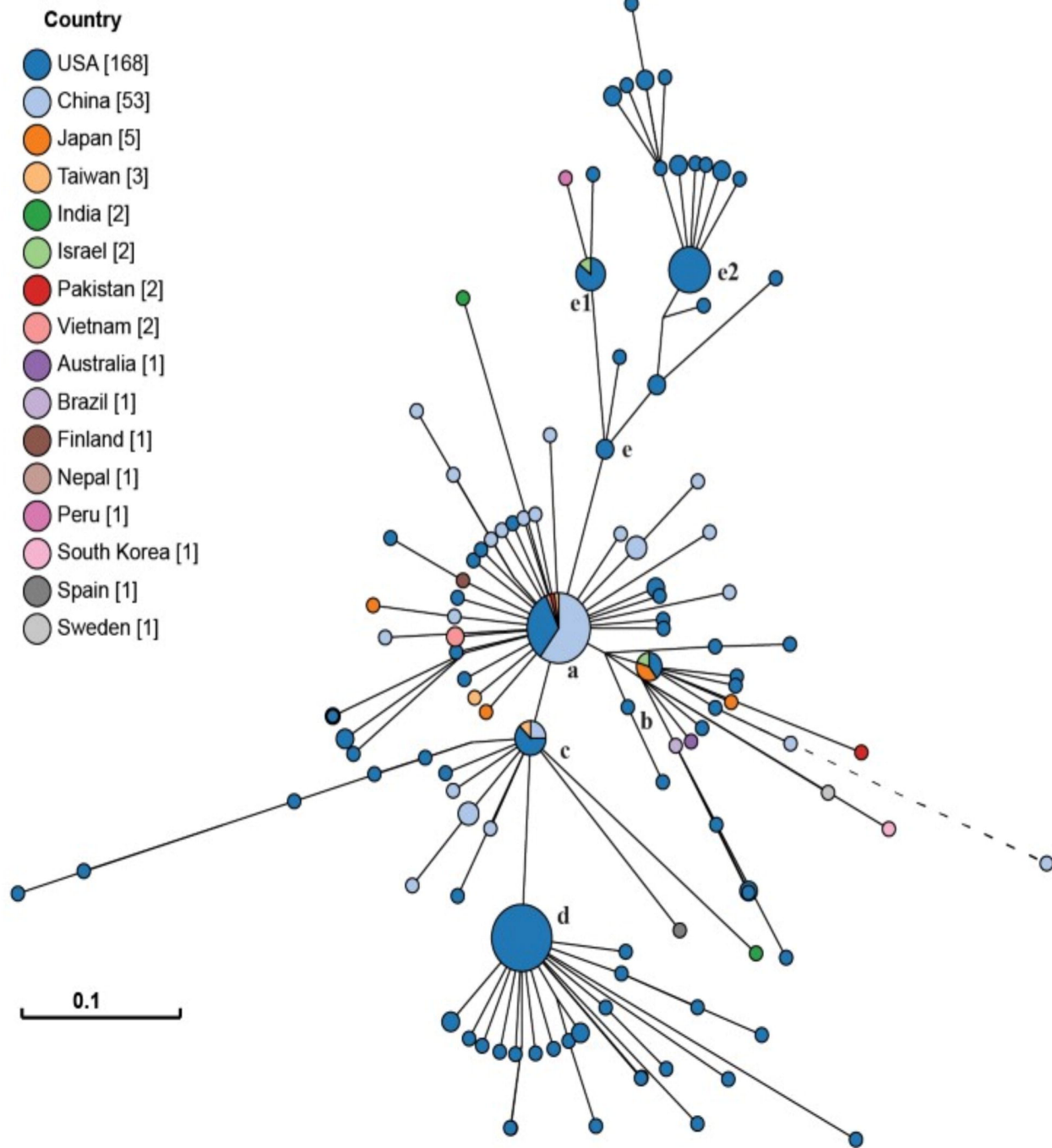
845

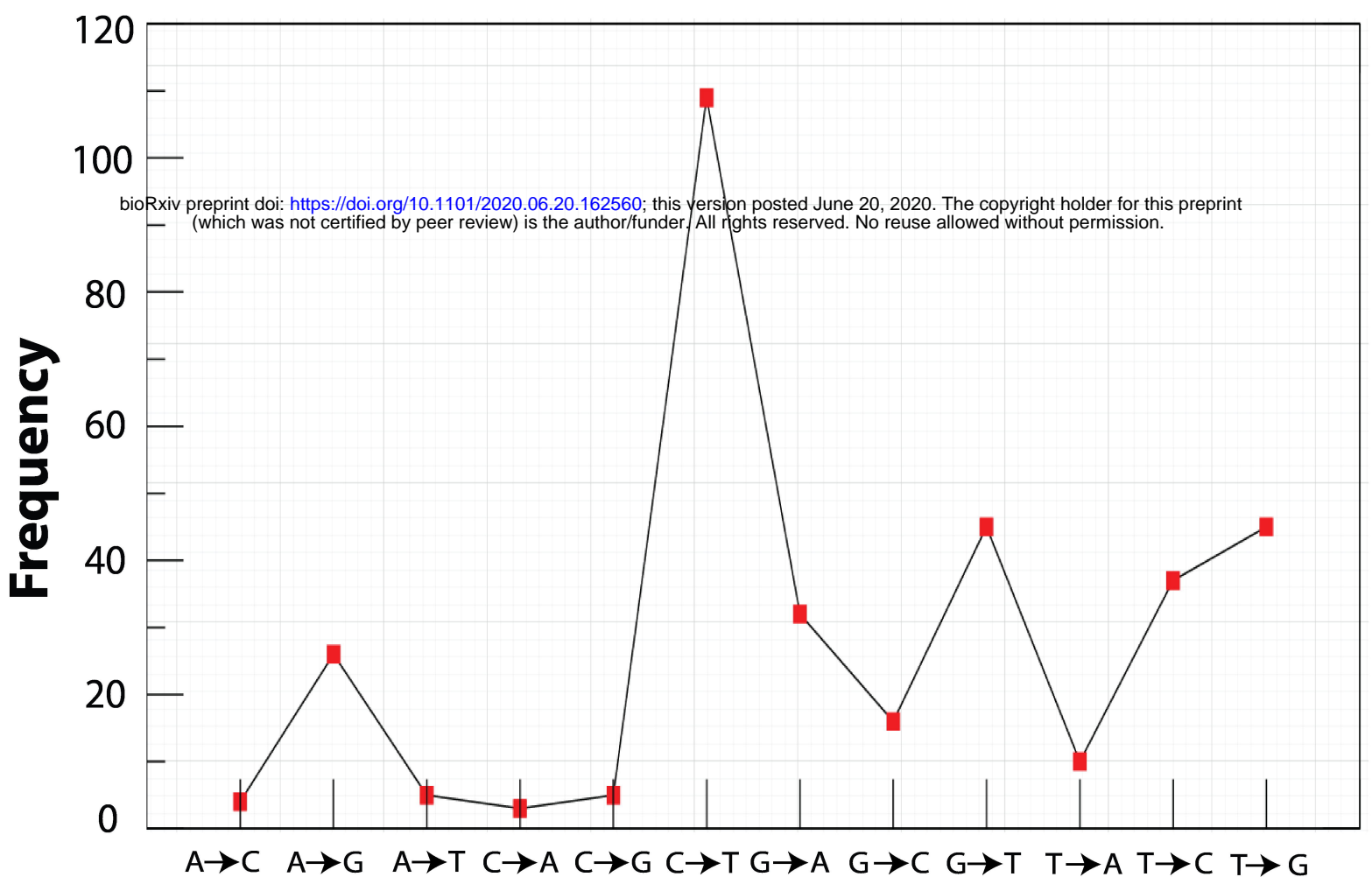
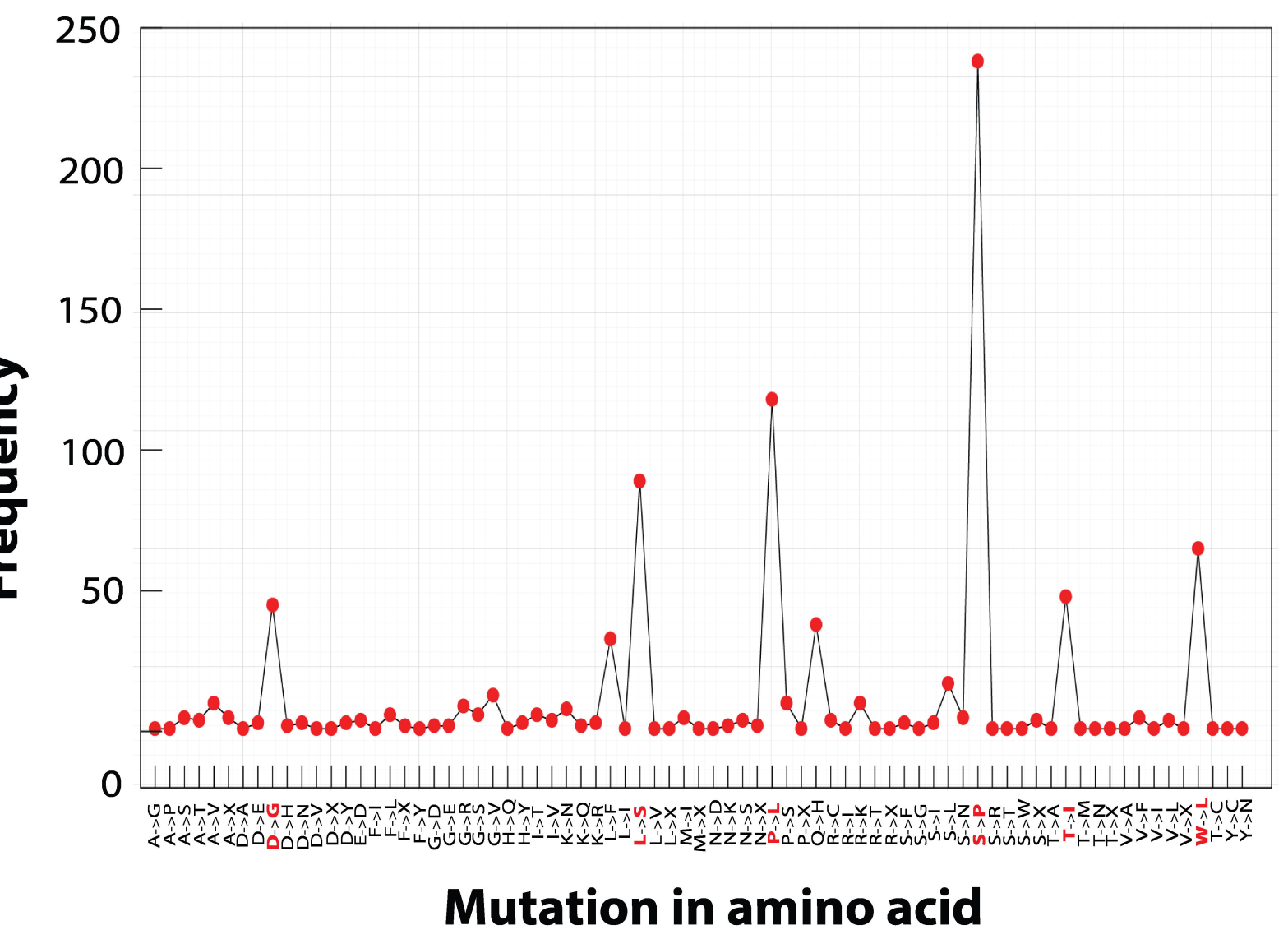
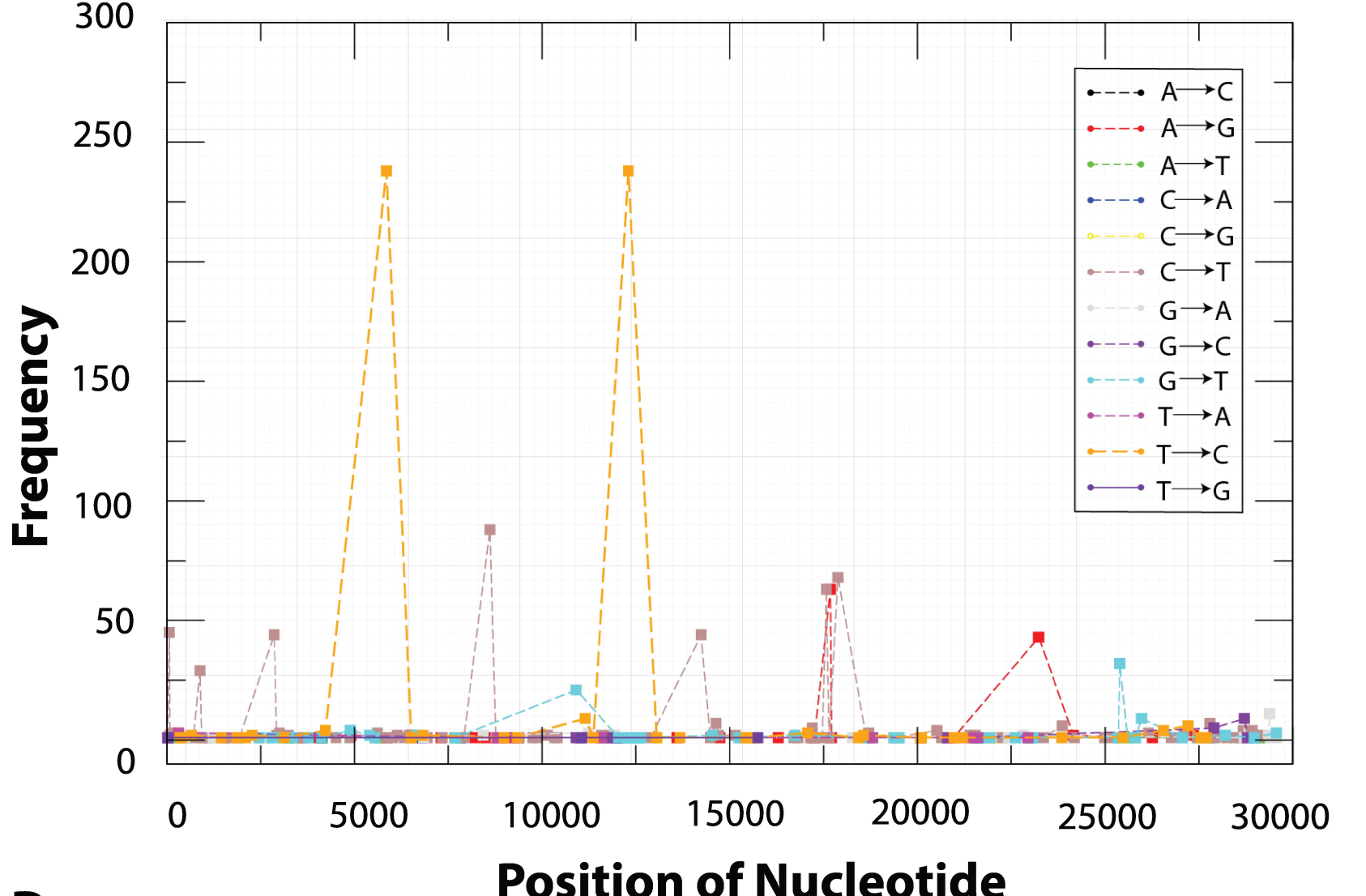
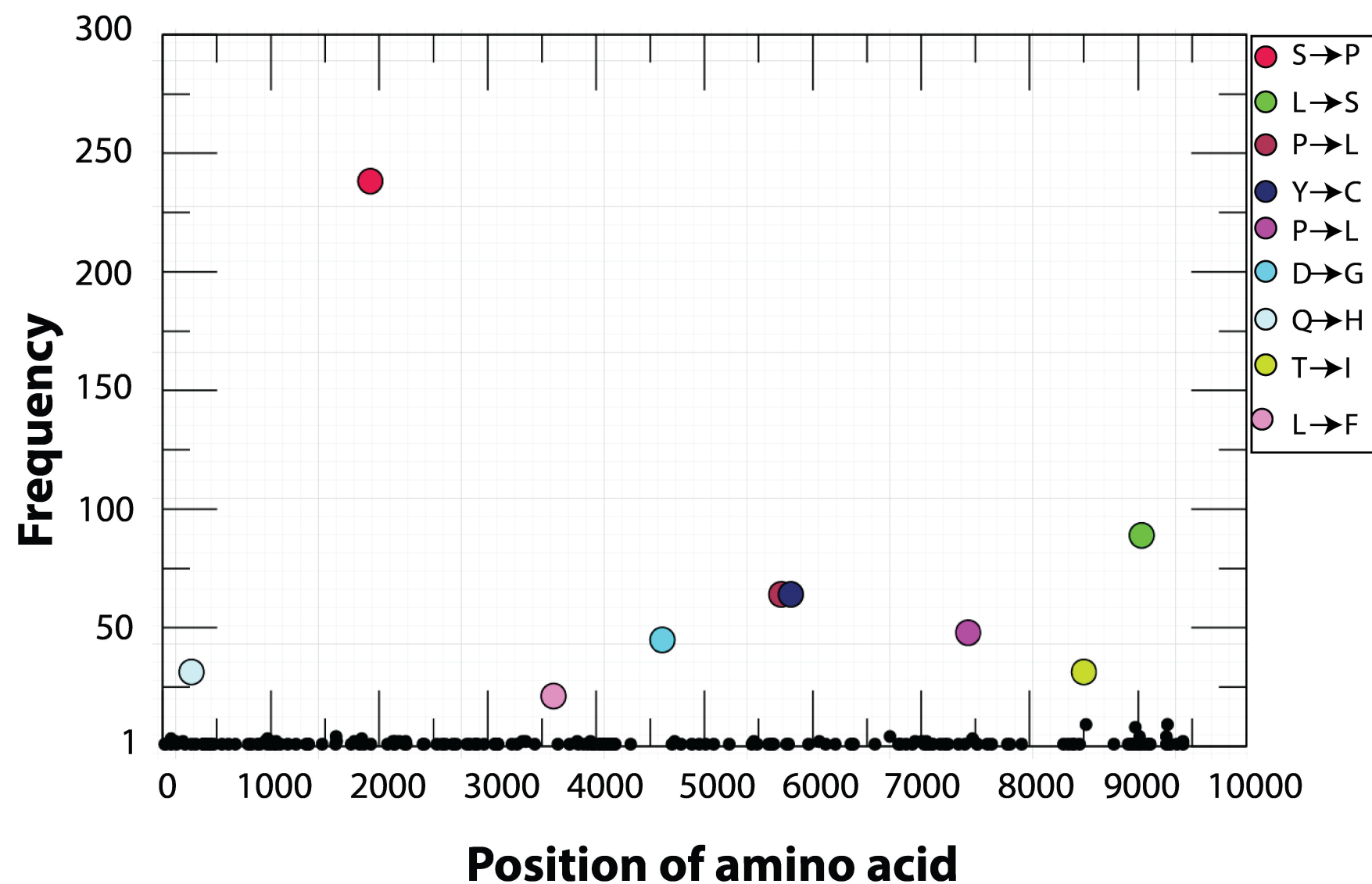
Tree scale: 0.01

A

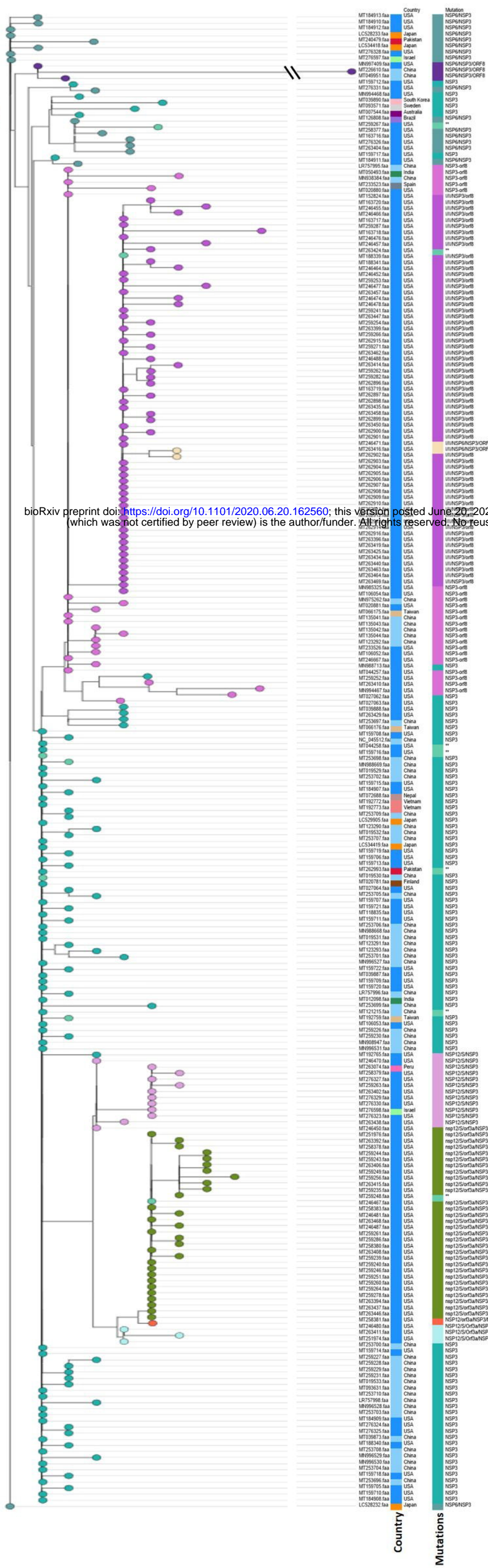


B.

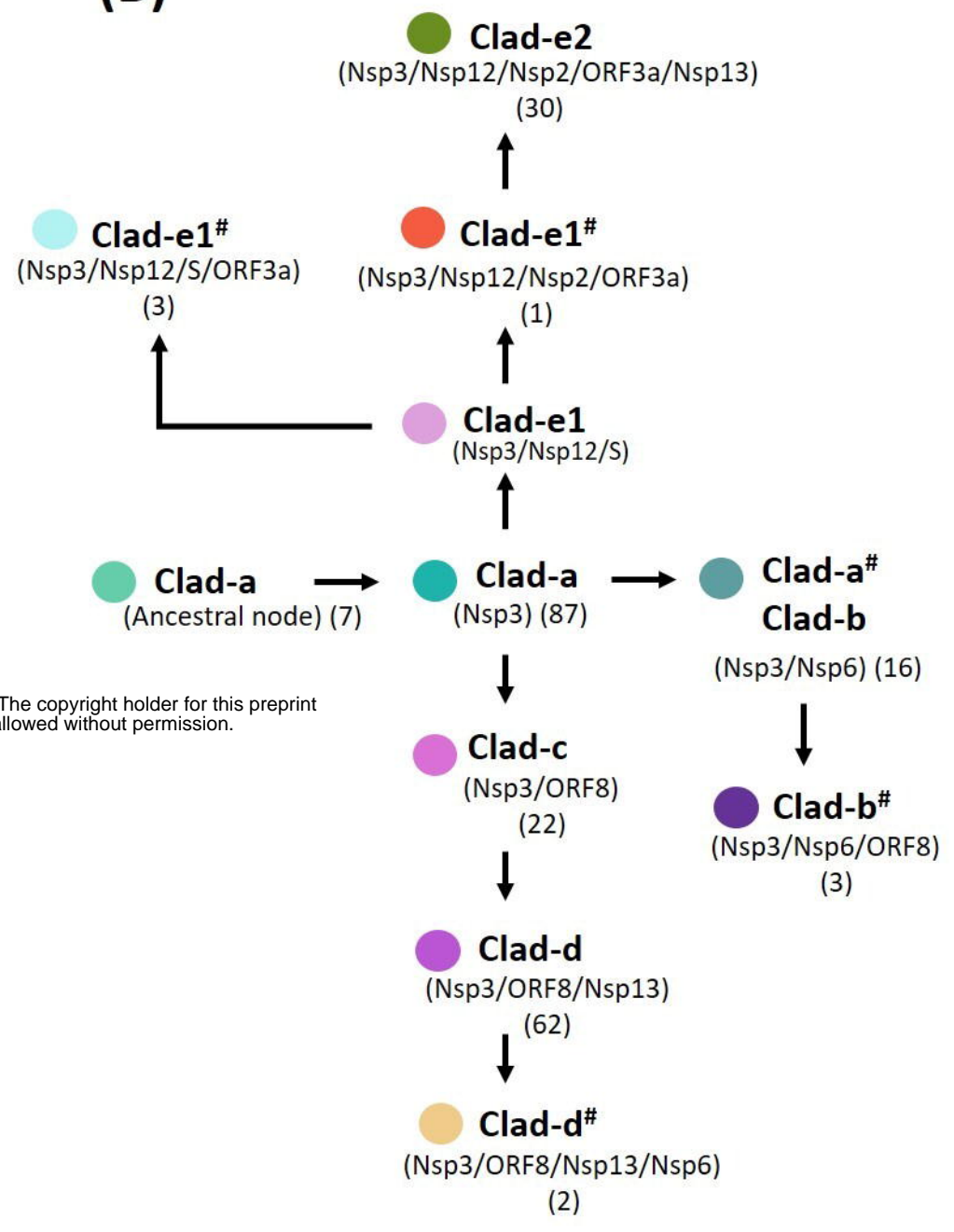


**A.****Single nucleotid polymorphism****C.****Mutation in amino acid****B.****D.**

(A)



(B)



#- Radiating daughter node

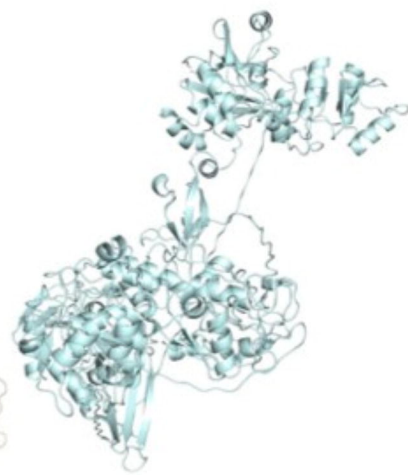
Country Legend

- USA
- China
- Australia
- Brazil
- India
- Japan
- Pakistan
- South Korea
- Finland
- Israel
- Taiwan
- Sweden
- Nepal
- Peru
- Spain
- Vietnam

bioRxiv preprint doi: <https://doi.org/10.1101/2020.06.20.162560>; this version posted June 20, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.



S-Protein



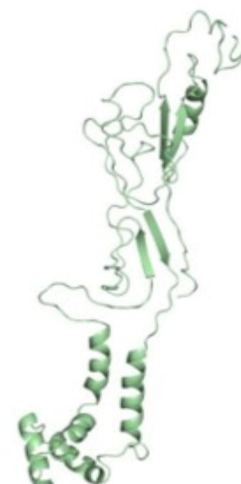
NSP3



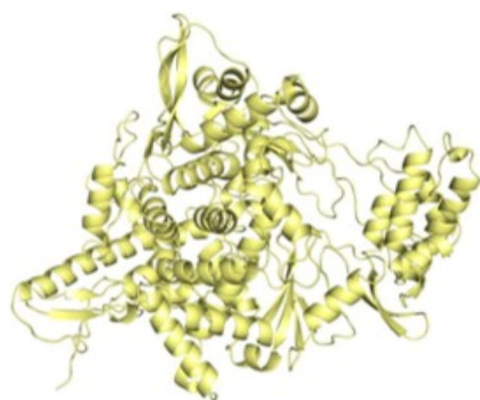
NSP6



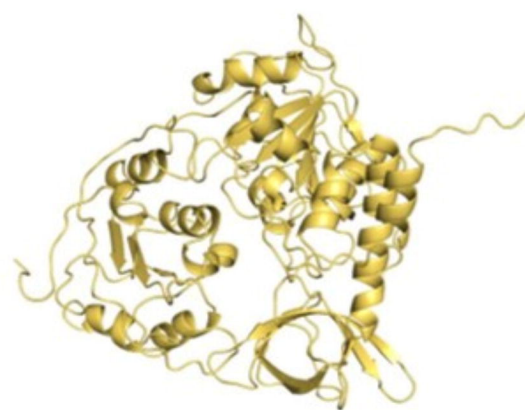
ORF8



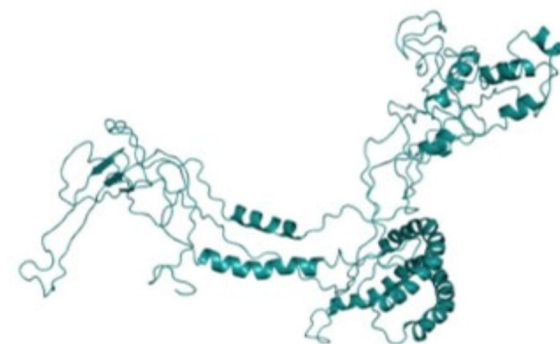
ORF3A



NSP12



NSP13



NSP2

