

1 **Full Title:**

2 **Molecular Evolution of SARS-CoV-2 structural genes:**  
3 **Evidence of positive selection in Spike Glycoprotein**

4

5 **Authors:** Xiao-Yong Zhan<sup>1</sup>, Ying Zhang<sup>1</sup>, Xuefu Zhou<sup>1</sup>, Ke Huang<sup>1</sup>, Yichao Qian<sup>1</sup>, Yang Leng<sup>1</sup>,

6 Leping Yan<sup>1</sup>, Bihui Huang<sup>1\*</sup>, Yulong He<sup>1\*</sup>

7

8 1 The Seventh Affiliated Hospital, Sun Yat-sen University, Shenzhen 517108, China

9

10 **Correspondence authors:**

11 \* Bihui Huang

12 ORCID: 0000-0003-3084-7096

13 Address: No.628, Zhenyuan Road, Guangming District, Shenzhen 518107, China

14 Tel: 86-755-81207035, Email: [huangbh7@mail.sysu.edu.cn](mailto:huangbh7@mail.sysu.edu.cn)

15 \*Yulong He

16 ORCID: 0000-0001-8930-8704

17 Address: No.628, Zhenyuan Road, Guangming District, Shenzhen 518107, China

18 Tel: 86-755-81201030, Email: [heyulong@mail.sysu.edu.cn](mailto:heyulong@mail.sysu.edu.cn)

## 19 **Abstract**

20 SARS-CoV-2 caused a global pandemic in early 2020 and has resulted in more than 8,000,000  
21 infections as well as 430,000 deaths in the world so far. Four structural proteins, envelope (E),  
22 membrane (M), nucleocapsid (N) and spike (S) glycoprotein, play a key role in controlling the entry  
23 into human cells and virion assembly of SARS-CoV-2. However, how these genes evolve during  
24 its human to human transmission is largely unknown. In this study, we screened and analyzed  
25 roughly 3090 SARS-CoV-2 isolates from GenBank database. The distribution of the four gene  
26 alleles is determined: 16 for E, 40 for M, 131 for N and 173 for S genes. Phylogenetic analysis shows  
27 that global SARS-CoV-2 isolates can be clustered into three to four major clades based on the  
28 protein sequences of these genes. Intra-genic recombination event isn't detected among different  
29 alleles. However, purifying selection has conducted on the evolution of these genes. By analyzing  
30 full genomic sequences of these alleles using codon-substitution models (M8, M3 and M2a) and  
31 likelihood ratio tests (LRTs) of codeML package, it reveals that codon 614 of S glycoprotein has  
32 subjected to strong positive selection pressure and a persistent D614G mutation is identified. The  
33 definitive positive selection of D614G mutation is further confirmed by internal fixed effects  
34 likelihood (IFEL) and Evolutionary Fingerprinting methods implemented in Hyphy package. In  
35 addition, another potential positive selection site at codon 5 in the signal sequence of the S protein  
36 is also identified. The allele containing D614G mutation has undergone significant expansion during  
37 SARS-CoV-2 global pandemic, implying a better adaptability of isolates with the mutation.  
38 However, L5F allele expansion is relatively restricted. The D614G mutation is located at the  
39 subdomain 2 (SD2) of C-terminal portion (CTP) of the S1 subunit. Protein structural modeling  
40 shows that the D614G mutation may cause the disruption of salt bridge among S protein monomers

41 increase their flexibility, and in turn promote receptor binding domain (RBD) opening, virus  
42 attachment and entry into host cells. Located at the signal sequence of S protein as it is, L5F mutation  
43 may facilitate the protein folding, assembly, and secretion of the virus. This is the first evidence of  
44 positive Darwinian selection in the *spike* gene of SARS-CoV-2, which contributes to a better  
45 understanding of the adaptive mechanism of this virus and help to provide insights for developing  
46 novel therapeutic approaches as well as effective vaccines by targeting on mutation sites.  
47

## 48 **Introduction**

49 Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), the causative agent of an  
50 emerging coronavirus disease (COVID-19) that has caused more than 430,000 deaths, is still a  
51 serious global pandemic currently. The genome of SARS-CoV-2 is consisting of a single-stranded  
52 and positive-sense RNA of around 30 kb in length with a 5' cap and 3'-polyA tail. It shows that  
53 SARS-CoV-2 genome possesses six major open reading frames (ORFs) that encodes 27 different  
54 proteins, in which four are structural proteins named Envelope (E), Membrane (M), Nucleocapsid  
55 (N) and Spike (S). Many studies have demonstrated important functions of these proteins in virus  
56 entry, transcription and virion particle assembly of SARS-CoV-2. The E protein is a small envelope  
57 protein with 75 amino acids. Given that a close genetic relationship between SARS-CoV-2 and  
58 SARS-CoV, functions of this protein may include virion assembly and morphogenesis[1]. In  
59 addition, induction of apoptosis of host cells might be another crucial function of SARS-CoV-2 E  
60 protein, thus making it a potential determinant of viral pathogenesis [2]. M protein, consisting of  
61 222 amino acids, is the most abundant component of the viral envelope and plays a key role in the  
62 virion assembly[3]. N protein, composed of 419 amino acids, may form complexes with genomic  
63 RNA, interact with the viral membrane protein, and play a critical role in enhancing the efficiency  
64 of virus transcription and assembly[4]. S protein, consisting of 1,273 amino acids, is the most  
65 important factor that mediates virus entry and a primary determinant of cell tropism and  
66 pathogenesis of SARS-CoV-2[5].

67 Many studies demonstrated SARS-CoV-2 underwent the evolution and some genetic evolutionary  
68 features have been reported[6]. The whole genomic sequence of SARS-CoV-2 has 79.6% identity  
69 with SARS-CoV and 96% with a bat SARS-related coronavirus (SARSr-CoV), RaTG13. Although

70 no positive time evolution signal was found between SARS-CoV-2 and RaTG13, the SARS-CoV-  
71 2 shows a strong positive temporal evolution relationship with bat-SL-CoVZC45, which has a  
72 slightly less identical genomic sequence (87.5%) than RaTG13 [7]. Combining the phylogenetic  
73 analysis of full-length genomes of coronaviruses, a potential bat origin of SARS-CoV2 is indicated  
74 [8]. A recent study reported that *spike* (S) gene (coding gene of S protein) of SARS-CoVs from  
75 their natural reservoir host, the Chinese horseshoe bat (*Rhinolophus sinicus*), has coevolved with *R.*  
76 *sinicus* angiotensin converting enzyme 2 (ACE2) via positive selection[9]. A single-stranded  
77 positive-sense RNA virus as it is, SARS-CoV-2 causes global pandemic within half a year,  
78 suggesting it may evolve rapidly. However, the evolution of SARS-CoV-2 based on structural genes  
79 from human to human transmission has not been investigated in detail. The primary purpose of this  
80 work is to study the evolutionary pattern of the four structural genes of SARS-CoV-2 derived from  
81 a global isolate collection including the E, M, N and S. Various molecular evolution and selection  
82 analysis approaches were employed to identify the phylogeny of the four structural proteins and  
83 potential selection effects on these genes. Hereby, our study reveals that intragenic recombination  
84 does not contribute to the evolution of these genes while purifying selection is the main evolutionary  
85 force. Moreover, a D614G mutation in the S protein is operated by strong positive selection and  
86 may be responsible for the quick spread of SARS-CoV-2 globally. Additionally, another potential  
87 L5F mutation may also be operated by positive selection, but with relatively less strong pressure as  
88 compared to D614G.

89

## 90 **Materials and Methods**

### 91 **SARS-CoV-2 isolates**

92 Complete full-length genomic sequences of SARS-CoV-2 were downloaded from 2019 Novel  
93 Coronavirus Resource (2019nCoV) in China National Center for Bioinformation. All of which  
94 were also uploaded to the NCBI GenBank database. The sequences were manually checked and  
95 finally a total of 3090 isolates were selected and verified for the present study. These isolates were  
96 collected from December 24, 2019 to April 24, 2020 in the different geographical locations  
97 including China, USA, Japan, Pakistan, Australia, Greece, German, Peru, Turkey, Kazakhstan, Iran,  
98 Serbia, Thailand, Nederland, Sri Lanka, Czech, Malaysia, India etc. Detailed information of these  
99 isolates including the GenBank accession number or biosample number is summarized in [S1 Table](#).

100

## 101 **Sequence analysis of the four structural genes and proteins**

102 The E, M, N, S gene sequences were extracted from SARS-CoV-2 global isolate collection and  
103 aligned by the MEGA X package using Muscle (codons) parameters [10]. Because some regions of  
104 genomic sequences of SARS-CoV-2 couldn't be exactly identified, in which nucleic acid bases are  
105 shown as degenerate bases (e.g. N, R, Y), we were unable to obtain all of the four structural gene  
106 sequences from an isolate sometimes. Allele type and DNA sequence polymorphism analyses were  
107 performed using DnaSP 6.12.03[11]. The protein sequences and polymorphism loci of these isolates  
108 were also aligned and analyzed with the MEGA X.

109

## 110 **Molecular evolution analysis**

111 An unrooted phylogenetic tree of the four structural proteins was constructed using the MEGA X  
112 package [10], and the evolutionary history was inferred using the Maximum Likelihood method,  
113 based on the JTT matrix-based model for E protein sequences, General Reversible Chloroplast +

114 Freq. model for M, JTT matrix-based model for N and Jones et al. w/freq. model for S protein  
115 sequences. Model selection was conducted in MEGA X. Bootstrap values were estimated by 1000  
116 replications. Initial tree(s) for the heuristic search were obtained automatically by applying  
117 Neighbor-Join and BioNJ algorithms to a matrix of pairwise distances estimated using each model  
118 mentioned above. The tree is drawn to scale, and FigTree V1.4 was utilized to form cladogram  
119 branches (<http://tree.bio.ed.ac.uk/software/figtree/>). The aligned DNA sequences were also  
120 screened using RDP4 software to detect intragenic recombination among the alleles of each  
121 structural gene[12]. Six methods implemented in the RDP4 were utilized. These methods are RDP  
122 [12], GENECONV[13], BootScan [14], MaxChi[15], Chimaera [16], and SiScan [17]. Common  
123 settings for all methods include considering sequences as linear and setting statistical significance  
124 at the  $P < 0.05$  with Bonferroni correction for multiple comparisons and requiring phylogenetic  
125 evidence and polishing of breakpoints. Potential recombination events (PREs) were considered as  
126 those identified by at least two methods. Reticulate network tree of alleles of the four structural  
127 genes of SARS-CoV-2 was also generated by Splitstree4 [18]. Phi test implemented in Splitstree4  
128 was used to define probable recombination events. Tajima's D, Fu and Li's D\* and F\* tests were  
129 employed to test the mutation neutrality hypothesis of the whole gene as previously described by  
130 our research group[19]. These analyses were carried out using DnaSP 6.12.03[11]. A statistical  
131 significance level with  $P < 0.05$  is acceptable. The false discovery rate and 1000 replications in a  
132 coalescent simulation were applied for correcting multiple comparisons. Non-neutrality evolution  
133 was considered when identified by at least two out of three tests. Nonsynonymous and synonymous  
134 mutations of the alleles of the four structural genes were also calculated using MEGA X package  
135 [10].

136

### 137 **Analysis of positive selection based on codon**

138 The selection pressure operating the four structural genes of SARS-CoV-2 was searched by using  
139 the Maximum Likelihood (ML) method. Analyses were performed using a visual tool of codeml  
140 program, named EasyCodeML algorithm with site model [20]. Three nested models (M3 vs. M0,  
141 M2a vs. M1a, and M8 vs. M7) were compared and likelihood ratio tests (LRTs) were applied to  
142 access a better fit of codes. Model fitting was also performed using multiple seed values for  $dN/dS$   
143 and assuming the F3x4 model of codon frequencies. Positive selection is inferred when individual  
144 site or codon with ratio of nonsynonymous to synonymous mutations ( $dN/dS$  ratios) is greater than  
145 one ( $\omega > 1$ ). When the LRT is significant ( $p < 0.05$ ), Bayes empirical Bayes (BEB) (M8 model) and  
146 Naive Empirical Bayes (NEB) methods (M3 and M2a model) are further employed to identify  
147 amino acid residues that likely evolve under positive selection based on a posterior probability  
148 threshold of 0.95. Results from M8 model were taken as the standard as Yang *et al.* reported. M3  
149 model was used for the frequency distribution of codon class analysis as Yang *et al.*  
150 recommended[21]. HyPhy package was used to validate the result obtained by ML method[22].

151

### 152 **Structural modeling of the protein with positive selection sites**

153 Three-dimensional structures of proteins with positive selection sites were modeled using SWISS-  
154 MODEL (<http://swissmodel.expasy.org>) according to the most fitted protein template. Model  
155 quality was evaluated by QMEAN while the structure of the model was visualized by using PyMoL  
156 [23].

157



## 158 **Results and Discussion**

### 159 **Characteristics of SARS-CoV-2 isolates, structural gene and protein** 160 **sequences**

161 The 3090 SARS-CoV-2 isolates harbor only 16 unique alleles of E and 40 alleles of M, but an  
162 abundant number of alleles of N and S genes, which contain 131 and 173, respectively. These alleles  
163 correspond to 10, 14, 88 and 99 different amino acid sequences of E, M, N, and S proteins,  
164 respectively. Protein sequence comparisons of WH01 isolate with SARSr-CoV, bat-SL-CoVZC45  
165 isolate show 100% (75/75) identity in E, 98.65% (219/222) identity in M, 94.27% (395/419) identity  
166 in N and 80.06% (1171/1273) in S proteins, respectively. These results imply a close kinship  
167 between SARS-CoV-2 and bat SARSr-CoV, especially on E and M proteins. On the other hand, it  
168 indicates an extreme conservation of E and M proteins and their functions among coronaviruses[24].  
169 Further analysis revealed that there are 14 single nucleotide polymorphisms (SNPs) of E gene, but  
170 only 5 single amino acid polymorphic (SAP) loci in the E protein. Similar result was observed on  
171 M gene and protein, with 37 SNPs and 9 SAPs. In contrast, 126 SNPs and 75 SAPs are detected on  
172 N gene and protein, respectively. S protein, the most important factor that mediates virus entry by  
173 receptor binding and membrane fusion and determines the infection ability of SARS-CoV-2 [25],  
174 harbors 155 SNPs on the alleles and 90 SAPs in the protein. Considering the size of nucleotides and  
175 amino acid residues, N gene has the maximum sequence variability with 10.02% (126/1257) SNPs  
176 and 17.90% (75/419) SAPs, respectively. However, S gene has most pairwise nucleotide differences  
177 among the four structural genes, indicating a more genetic diversity of S gene (Table 1). A key  
178 player in the virus transcription and assembly as N protein is [26, 27], high sequence variability of  
179 the N protein may indicate a vast adaption of the virus during host transmission. Previous study

180 shows that high genetic variance has been found among bat SARS-CoVs, particularly in the S  
181 gene[9]. Similar, higher nucleotide diversity ( $\pi$ , a major parameter to define genetic diversity) of S  
182 gene is also detected on SARS-CoV-2 isolates, suggesting this may benefit virus survival in the host  
183 of human beings.

184

185 **Table 1. Summary of genetic diversity of the 4 structural genes of the SARS-CoV-2 isolates**

Gene	Sequence, n*	Sequence length	$h$	$\pi$	$S$	$\theta$	$\eta$
E	2928	228	16	0.00012	14	0.00475	15
M	2891	669	40	0.00018	37	0.00665	40
N	2253	1260	131	0.00056	126	0.01081	130
S	2339	3825	173	0.00075	155	0.00753	169

186

187  $h$ , Haplotypes,

188  $\pi$ , Nucleotide diversity

189  $S$ , Polymorphic sites

190  $\theta$ , Theta (per site) from  $S$ , population mutation ration

191  $\eta$ , Total number of mutations

192 \* Some bases of SARS-CoV-2 genomic sequences are not exactly identified; thus, the number of  
193 gene sequences were less than 3090.

194

## 195 **Distinct phylogenetic patterns of the four structural genes**

196 The phylogenetic analysis revealed that all SARS-CoV-2 E proteins form three clusters. Similar to

197 E protein, phylogenetic tree of SARS-CoV-2 M proteins is formed by three clusters with few

198 branches (Figs 1A and 1B). The results suggest both E and M genes may display a relatively high

199 conservation during coronavirus evolution. In contrast, SARS-CoV-2 N and S proteins show distinct

200 phylogenetic pattern as compared with that of E and M. Four and three main phylogenetic clusters

201 with various branches are identified in the N and S proteins, respectively (Figs 1C and D). Given

202 the crucial roles of N and S proteins in virus transcription, assembly, and entry to host cells, whether

203 SARS-CoV-2 isolates harbor different N and S variants (such as those clustered into different clades)

204 may influence their infection efficiency remains unknown, and requires further study.

205 **Purifying selection drives the evolution at whole structural gene levels**  
206 **of SARS-CoV-2 during its human to human transmission**

207 Although many studies demonstrated that recombination plays an important role on the emergence

208 of SARS-CoV-2 and its contribution to admit SARS-CoV-2 as a human infectious pathogen [28-

209 30], how this virus evolves during its global transmission has not been profiled yet. Therefore, we

210 first analyzed intragenic recombination events of each structural gene using RDP4. The results

211 indicate there were no recombination events occurred among the alleles of each gene (data not

212 shown). Recombination event were also assessed through reticulate network tree by phi test in

213 SplitsTree4. Although some internal nodes are noticed in N and S alleles, no significant evidence

214 for recombination is validated of each gene by Phi test ( $p > 0.05$ ) (Fig 2). It indicates a relative stable

215 state of SARS-CoV-2 during its transmission although a possible genetic interaction of different

216 isolates might have occurred when it became a global pandemic [31, 32]. In addition, Tajima's D,

217 Fu and Li's D\* and F\* statistics were calculated to examine the mutation neutrality hypothesis of

218 the four structural genes of SARS-CoV-2. The results reveal that the evolution of all four genes

219 does not match the neutral hypothesis, but favor purifying selection (Table 2 and Fig 3). The average

220 of all pairwise  $dN/dS$  ratios ( $\omega$ ) among the alleles of each structural gene of SARS-CoV-2 is 0.5443

221 in E, 0.1562 in M, 0.07978 in N, and 0.4980 in S gene, respectively. All together, these results

222 suggest that at the whole gene level, inconsistent purifying selection is the main evolution force

223 (Table 2). Li et al studied the origin of SARS-CoV-2 and showed evidence of strong purifying

224 selection in the S and other genes among bat, pangolin and human coronaviruses, indicating similar

225 strong evolutionary constraints in different host species [33]. Similarly, our results suggest purifying  
226 selection drives the evolution at the whole structural gene level of SARS-CoV-2 during its  
227 transmission from human to human. This result also implies that in general, the genetic variation on  
228 these structural genes will not confer a significant disadvantage on the virus survival, and ratios  
229 reflect general variability of these genes and proteins. Considering that no recombination happened,  
230 nonsynonymous mutations would be removed at a great rate during the virus transmission [34].

231

232 **Table 2. Summary of neutrality for the four structural genes in SARS-CoV-2 isolates**

Gene	Tajima'sD	Fu and Li's D* test t	Fu and Li's F * test	dN	dS	dN/dS ( $\omega$ )	Selection
E	-2.29974, P<0.01	-3.18477, P<0.02	-3.38505, P<0.02	0.006836	0.1256	0.5443	Purifying selection
M	-2.74611, P<0.001	-5.64276, P<0.02	-5.50855, P<0.02	0.001294	0.008296	0.1562	Purifying selection
N	-2.87598, P<0.001	-9.67153, P<0.02	-7.95879, P<0.02	0.000251	0.003146	0.07978	Purifying selection
S	-2.87646, P<0.001	-11.01171, P<0.02	-8.59037, P<0.02	0.000609	0.001223	0.4980	Purifying selection

233

234 **SARS-CoV2 S gene is operated by positive selection at a definitive**  
235 **codon located at the C-terminal portion of S1 subunit and a potential**  
236 **codon located at the signal sequence**

237 Guo et al. reported that the S gene of SARS-CoV populations in their natural host, Chinese  
238 horseshoe bat (*Rhinolophus sinicus*), has evolved through positive selection at some codons[9]. As  
239 mentioned above, at the whole gene level, purifying selection is the main force driving the evolution  
240 of studied genes. Whether positive selection pressure accelerates the diversification of the structural  
241 genes of SARS-CoV-2 remains unclear. Therefore, we used codon-substitution models to estimate  
242 the ratio of nonsynonymous over synonymous substitutions ( $dN/dS$ ), also known as  $\omega$ . The role of

243 recombination in the polymorphism of four genes is excluded because no intragenic recombination  
244 was detected (Fig 2). By using ML model, we don't find any codon of E and M gene subjecting to  
245 positive selection obviously (data not shown). However, a potential positive selection site 208A in  
246 N gene is identified by using M3 model, but not by any other models especially the M8 model,  
247 suggesting a limited amount of evidence of positive selection in N gene (S1 Table). For the S gene,  
248 we found the average  $\omega$  is 0.37199 calculated by M0 model of the codeML package, suggesting that  
249 purifying selection was a major force operating the evolution of the S gene during its transmission  
250 among human beings. In three LRTs, all alternative models (M3, M2a, M8) are significantly better  
251 fit ( $P < 10^{-4}$ ) than relevant null models (M0, M1a, M7), indicating that some sites of S were subjected  
252 to strong positive selection ( $\omega = 18.22175 - 20.61283$ ) (Table 3). A single positive selection site (614D)  
253 is identified in the S gene with posterior probability of 1.000 in all the three models [21], a clear  
254 evidence showing that this site is still experiencing positive selection when the virus transmitted  
255 from human to human. The result is also validated using internal fixed effects likelihood (IFEL) and  
256 Evolutionary Fingerprinting methods implemented in HyPhy package (Fig 4) [35-37]. To our  
257 surprise, the positive selection site is not located at the receptor binding domain (RBD) or receptor  
258 binding motif (RBM) as we anticipated, which play the most important role in virus-receptor  
259 interaction and virus entry into host cells [38]. This result suggests that a relatively genetic stability  
260 of this motif would benefit the virus survival. Intriguingly, the site under positive selection pressure  
261 always has a D614G (for the S gene is 1841A>G) mutation, implying such mutation may enhance  
262 virus adaptability in human hosts. Another potential positive selection site at codon 5 is also  
263 identified, and a L5F mutation (for the S gene is 13C>T) is always found, with posterior  
264 probabilities greater than 0.95, 0.93 and 0.92 (critical values) calculated by M3, M2a and M8 models

265 (Table 3), respectively. Similar result was also confirmed by Evolutionary Fingerprinting method  
266 (S1 Fig). Considering signal sequence (SS) is a short hydrophobic peptide that plays an important  
267 role in guiding viral protein into the endoplasmic reticulum (ER) for proper folding and assembly  
268 [39], we postulate that L5F mutation may increase hydrophobicity of the SS, thus facilitating the  
269 entry of S protein into ER for folding and assembly, and in turn secretion of the virus.  
270

271 **Table 3. Log-likelihood values and parameter estimates for the SARS-CoV-2 S gene sequences**

Model	Ln L	Estimates of parameters	Model compared	LRT P-value	Positive sites
		p0=0.96797, p1=0.02883, p2=0.00320			<b>5 L 0.958*</b> ,28 Y 0.850,221 S 0.901, <b>614 D</b>
M3 (discrete)	-6766.339162	$\omega_0=0.26126$ , $\omega_1=2.70530$ , $\omega_2=$ <b>20.61283</b>			<b>1.000***</b> ,677 Q 0.891
M0 (one ratio)	-6790.072925	$\omega_0=0.37199$	M0 vs. M3	0.000000001	Not Allowed
		p0=0.81731, p1=0.17872, p2=0.00397			5 L 0.9258,28 Y 0.812,221 S 0.832, <b>614 D</b>
M2a(selection)	-6766.432802	$\omega_0=0.17504$ , $\omega_1=1.00000$ , $\omega_2=$ <b>18.76936</b>			<b>1.000***</b> ,677 Q 0.828
		p0=0.70461, p1=0.29539			
M1a (neutral)	-6778.770190	$\omega_0=0.04395$ , $\omega_1=1.00000$	M1a vs. M2a	0.000004385	Not Allowed
		p0=0.99578, p=0.40368, q=0.82224			5 L 0.931,28 Y 0.817,221 S 0.831, <b>614 D</b>
M8(beta& $\omega$ )	-6768.829411	p1= 0.00422, $\omega=$ <b>18.22175</b>			<b>1.000***</b> ,677 Q 0.828
M7(beta)	-6779.230494	p=0.00857, q=0.02623	M7 vs.M8	0.000030400	Not Allowed

272

273 LnL is the log likelihood;  $\omega$  is ratio of  $dN/dS$ , LRT P-value indicates the value of chi-square test; Parameters indicating positive selection are presented in bold;

274 Positive selection sites were identified by the Bayes empirical Bayes (BEB) methods under M8 model. The posterior probabilities  $(p) \geq 0.80$  are shown,  $(p) \geq 0.95$

275  $(p) \geq 0.99$ , and  $(p) = 1.000$  are indicated by \*, \*\* and \*\*\*, respectively. Yang *et al.* recommended that results from M8 model were preferred to find sites under positive

276 selection pressure.

## 277 **Evolutionary relationship of S gene alleles with or without D614G and** 278 **L5F mutation**

279 Phylogenetic tree of S gene alleles was derived to test the evolutionary relationship among the alleles  
280 with or without D614G mutation. As shown in **Fig 5A**, the 173 alleles of the S gene could be  
281 clustered into four clades. Alleles with D614G mutation could be found in all 4 clades, among which  
282 a dominant one contains 79 out of 85 alleles with such mutation. The remaining 6 mutated S alleles  
283 are distributed in other 3 clades. The result suggests a potential common ancestor for the majority  
284 of S alleles with D614G mutation, while some other maybe derived from alternative ancestors. This  
285 result is also supported by the parsimony network of S gene alleles using PopART  
286 (<http://popart.otago.ac.nz>) [40]. Two central alleles (representative virus isolates are WH01 and  
287 GZMU0019) and associated alleles around them form a star scattering network, suggesting that the  
288 S gene may have two potential origins (**Fig 5B**). All S alleles with D614G mutation are closely  
289 related (with a few point mutations), and comprise a scattered star structure, suggesting the  
290 expansion of SARS-CoV-2 population with D614G mutation on S gene. In contrast, alleles of the  
291 N gene show a single ancestor analyzed by parsimony network though 3 phylogenetic clades are  
292 identified (**S2 Fig**).

293 A total of 5 alleles with L5F mutation are found and all of them are in one clade, accounting for  
294 83.33% of all alleles in the clade (**S3A Fig**). Further parsimony network analysis reveals that S  
295 alleles with L5F mutation are not closely related, but distribute in both WH01 and GZMU0019  
296 haplotype groups (**S3B Fig**). No scattered star structure of these alleles can be formed, indicating  
297 L5F mutation might arise from independent origins other than that of D614G mutants. Limited



298 number of alleles with L5F mutation identified so far also suggests that L5F might subject to  
299 relatively less strength of the pressure and is still at early stage of positive selection.

300

### 301 **Frequency of S allele with D614G mutation increased in SARS-CoV-2** 302 **isolates during human to human transmission**

303 Considering that mutation of a positive selection site should be beneficial to the survival of the  
304 individuals carrying the mutation, we postulate that the D614G (1841A>G) mutation may help the  
305 spread of SARS-CoV-2. Some evidence has been obtained from the haplotype network of S alleles  
306 mentioned above (Fig 5B). S gene haplotypes (alleles) with D614G mutation (representative isolate  
307 GZMU0019) have evolved many subtypes and comprise a star structure with GZMU0019 in the  
308 center. This starburst pattern with one haplotype in the center and many other haplotypes  
309 surrounding the central haplotype suggests a signature of rapid population expansion [41]. To  
310 further study whether SARS-CoV-2 isolates with D614G mutation have advantage in survival  
311 during its transmission among human beings, we calculated the frequencies of S alleles carrying  
312 D614G mutation in each week from the collected SARS-CoV-2 isolates from December 24, 2019  
313 to April 20, 2020 (17 weeks). Detailed information of these isolates including collection date,  
314 collection region and accession or biosample numbers is summarized on [S3 and S4 Tables](#).

315 In 173 S gene alleles, 85 carry D614G mutation, accounting for 49.13% of all. Similarly, 47 out 99  
316 S proteins carry D614G mutation, accounting for 47.47% of all. The first two isolates,  
317 GWHABKF00000001 and WH01 (isolated in December 24, 2019 and December 26, 2019,  
318 respectively), carry 614D in the S protein, while the first SARS-CoV-2 isolate with a D614G  
319 mutation is GZMU0019 in our collected dataset, isolated from a patient with COVID-19 in

320 Guangzhou, Guangdong Province of China on February 5, 2020 (week 7 in our dataset). After that,  
321 except for week 9 and week 10 (possibly due to the small number of samples and sampling  
322 deviation), a spread trend that more and more proportion of isolates carry the D614G mutation in  
323 the S protein stands out. In the week 17, the last week of our dataset, 91.11% of SARS-CoV-2  
324 isolates carry this mutation (S3 Table, Fig 6A). Further analysis reveals that the frequency of D614G  
325 mutation in the S gene was steadily increasing when combining data from week 6 to 17 (S3 Table,  
326 Fig 6B). To exclude the influence of sample size on the result (in some weeks, only 4-6 isolates  
327 were collected in the dataset), we reorganized the dataset by taking both the sample size and  
328 sampling time into account. Various panels of 200-300 isolates were studied and similar results  
329 were observed (S4 Table, Figs 6C and D). Taken together, these results suggest that SARS-CoV-2  
330 isolates with D614G mutation may increase their ability to transmit, and contribute to the rapid  
331 spread of this virus to the world.

332

### 333 **D614G mutation of S gene may destabilize S protein trimer and** 334 **promote receptor binding and membrane fusion**

335 The positive selected D614G mutation might play an important role for the adaptability of SARS-  
336 CoV-2 in both the host and the virus population[42]. Another explanation is that the mutation is  
337 driven by specific interaction between high level of virus sequence divergence and polymorphic  
338 host receptors or interacting proteins[43]. S protein is the key determinant for the tissue tropism and  
339 host range and specificity of coronavirus such as SARS-CoV-2. The virus infects host cells through  
340 the interaction between the S protein and its cellular receptor, named ACE2 [8]. In this process,  
341 virus entry requires the precursor S protein cleaved by cellular proteases including trypsin, furin,

342 transmembrane serine protease 2 (TMPRSS2), or endosomal cathepsin L, which generate the  
343 receptor binding subunit S1 and the membrane fusion S2 [44-46]. From structural studies in both  
344 SARS-CoV and SARS-CoV-2, receptor binding domain (RBD) located at the C-terminal of S1 and  
345 the adjacent N-terminal domain (NTD) are relatively flexible, which is the feature required for  
346 receptor recognition and subsequent membrane fusion[47, 48]. We found that the D614G mutation  
347 is located at the subdomain 2 (SD2) that at the C-terminal of RBD and close to the two potential  
348 cleavage sites between S1 and S2 [48] (Fig 7A). Considering that positive selection is usually  
349 beneficial to the survival of the individual carrying the mutation, we speculate that the D614G  
350 mutation may facilitate structural conformation change to promote receptor binding or membrane  
351 fusion[5, 44], and in turn improving the infection efficiency. From the latest cryo-electron  
352 microscopy (cryo-EM) structure of SARS-CoV-2 S protein, the negatively charged sidechain of  
353 D614 points towards the positively charged sidechain of K854 from the neighboring monomer (Fig  
354 7B) [48]. The distance between the closest atoms of the two residues is 2.6 Å, which is an optimal  
355 distance to form salt bridge (Fig 7C). From the modelled structure with D614G mutation, the  
356 distance is increased to 5.2 Å (Fig 7D), which would potentially abolish the salt bridge and  
357 destabilize the integrity of the S trimer in wild type. It has been reported that human receptor ACE2  
358 binds to an “open” conformation of S protein, where RBD move away from the core structure and  
359 expose its receptor binding surface. The entire S trimer then undergoes a serial of dramatic  
360 conformation changes, including cleavages between S1 and S2, disassociation of S1 and post-fusion  
361 transformation of S2 [49, 50]. Changes including mutations at cleavage sites and adding internal  
362 crosslinks in S trimer would keep the protein in a stable and “closed” conformation where the  
363 receptor binding surface of RBD is inaccessible [48, 51]. Therefore, we hypothesize that the

364 highly transmissible D614G mutation driven by the positive selection through evolution promotes  
365 accessibility of RBD by losing a critical salt bridge between the S protein monomers, which  
366 subsequently triggers membrane fusion upon ACE2 binding.

367

## 368 **Conclusions**

369 We present modern molecular evolution analyses on a large and comparative set of SARS-CoV-2  
370 structural gene sequences, derived from an international collection of SARS-CoV-2 isolates.

371 Distinct phylogenetic patterns of four structural proteins of SARS-CoV-2 are depicted. Protein  
372 sequence comparisons show E and M genes exhibit a relatively close relationship to bat SARSr-  
373 CoV, suggesting the evolution conservation of these two genes. In contrast, relatively high genetic  
374 variation is observed in N and S proteins among SARS-CoV-2 isolates, implying extensive  
375 adaptability of N and S genes. No clear intragenic recombination is detected of these four genes,  
376 suggesting that it is not the major force to drive the evolution of the four genes. However, our

377 analyses show purifying selection pressure may be the main force operating the evolution at whole  
378 gene levels of SARS-CoV-2 during its human to human transmission. We also identify a codon in

379 S gene definitively experiencing positive selection pressure, and always leads to the D614G  
380 mutation in S proteins. S alleles with D614G mutation have expanded rapidly among SARS-CoV-

381 2 isolates. D614G mutation significantly extends the distance between monomers in the S protein  
382 trimer, which may disrupt the salt bridge formed by D614 and K854 between monomers, promote

383 RBD opening, and facilitate the entry of the virus into host cells, thus contributing to the diffusion  
384 of this mutated alleles. Codon 5 of S gene is another potential positive selection site. Although a

385 limited number of alleles with L5F mutation is identified, it may potentially affect the assembly and

386 secretion of SARS-CoV-2. A close eye on L5F mutation may be required in case another expansion  
387 occurs. As S protein is a key target for SARS-CoV-2 vaccines, therapeutic antibodies, and  
388 diagnostics, the D614G mutation of S should be paid more attention. Owing that the exact  
389 mechanism remains unclear, further study should focus on the exact function of these mutation sites  
390 and how they affect the expansion of these mutated alleles on SARS-CoV-2.

391

## 392 **Acknowledgements**

393 This research was supported by National Natural Science Foundation of China (grant number  
394 31870001) to X.Y.Z.

395 **Conflict of Interest** The authors have declared no conflict of interests.

## 396 **References**

- 397 1. Liu DX, Yuan Q, Liao Y. Coronavirus envelope protein: a small membrane protein with  
398 multiple functions. *Cellular and molecular life sciences* : CMLS. 2007;64(16):2043-8. Epub  
399 2007/05/29. doi: 10.1007/s00018-007-7103-1. PubMed PMID: 17530462; PubMed Central  
400 PMCID: PMCPMC7079843.
- 401 2. Jimenez-Guardeno JM, Nieto-Torres JL, DeDiego ML, Regla-Nava JA, Fernandez-  
402 Delgado R, Castano-Rodriguez C, et al. The PDZ-binding motif of severe acute respiratory  
403 syndrome coronavirus envelope protein is a determinant of viral pathogenesis. *PLoS pathogens*.  
404 2014;10(8):e1004320. Epub 2014/08/15. doi: 10.1371/journal.ppat.1004320. PubMed PMID:  
405 25122212; PubMed Central PMCID: PMCPMC4133396.
- 406 3. Arndt AL, Larson BJ, Hogue BG. A conserved domain in the coronavirus membrane  
407 protein tail is important for virus assembly. *Journal of virology*. 2010;84(21):11418-28. Epub  
408 2010/08/20. doi: 10.1128/JVI.01131-10. PubMed PMID: 20719948; PubMed Central PMCID:  
409 PMCPMC2953170.
- 410 4. McBride R, van Zyl M, Fielding BC. The coronavirus nucleocapsid is a multifunctional  
411 protein. *Viruses*. 2014;6(8):2991-3018. Epub 2014/08/12. doi: 10.3390/v6082991. PubMed  
412 PMID: 25105276; PubMed Central PMCID: PMCPMC4147684.
- 413 5. Belouzard S, Millet JK, Licitra BN, Whittaker GR. Mechanisms of coronavirus cell entry  
414 mediated by the viral spike protein. *Viruses*. 2012;4(6):1011-33. Epub 2012/07/21. doi:  
415 10.3390/v4061011. PubMed PMID: 22816037; PubMed Central PMCID: PMCPMC3397359.
- 416 6. Wu F, Zhao S, Yu B, Chen YM, Wang W, Song ZG, et al. A new coronavirus associated  
417 with human respiratory disease in China. *Nature*. 2020;579(7798):265-9. Epub 2020/02/06. doi:

- 418 10.1038/s41586-020-2008-3. PubMed PMID: 32015508; PubMed Central PMCID:  
419 PMCPMC7094943.
- 420 7. Y. Z, S. Z, J. C, C. W, W. Z, B. Z. Analysis of variation and evolution of SARS-CoV-2  
421 genome. Journal of Southern Medical University. 2020;02:152-8. doi: 10.12122/j.issn.1673-  
422 4254.2020.02.23.
- 423 8. Zhou P, Yang XL, Wang XG, Hu B, Zhang L, Zhang W, et al. A pneumonia outbreak  
424 associated with a new coronavirus of probable bat origin. Nature. 2020;579(7798):270-3. Epub  
425 2020/02/06. doi: 10.1038/s41586-020-2012-7. PubMed PMID: 32015507; PubMed Central  
426 PMCID: PMCPMC7095418.
- 427 9. Guo H, Hu B-J, Yang X-L, Zeng L-P, Li B, Ouyang S-Y, et al. Evolutionary arms race  
428 between virus and host drives genetic diversity in bat SARS related coronavirus spike genes.  
429 2020:2020.05.13.093658. doi: 10.1101/2020.05.13.093658. bioRxiv.
- 430 10. Kumar S, Stecher G, Li M, Knyaz C, Tamura K. MEGA X: Molecular Evolutionary Genetics  
431 Analysis across Computing Platforms. Molecular biology and evolution. 2018;35(6):1547-9.  
432 Epub 2018/05/04. doi: 10.1093/molbev/msy096. PubMed PMID: 29722887; PubMed Central  
433 PMCID: PMCPMC5967553.
- 434 11. Rozas J, Ferrer-Mata A, Sánchez-DelBarrio JC, Guirao-Rico S, Librado P, Ramos-Onsins  
435 SE, et al. DnaSP 6: DNA Sequence Polymorphism Analysis of Large Datasets. 2017;34(12).
- 436 12. Martin DP, Murrell B, Khoosal A, Muhire B. Detecting and Analyzing Genetic  
437 Recombination Using RDP4. Methods in molecular biology. 2017;1525:433-60. doi:  
438 10.1007/978-1-4939-6622-6\_17. PubMed PMID: 27896731.
- 439 13. Padidam M, Sawyer S, Fauquet CM. Possible emergence of new geminiviruses by

- 440 frequent recombination. *Virology*. 1999;265(2):218-25. Epub 1999/12/22. doi:  
441 10.1006/viro.1999.0056. PubMed PMID: 10600594.
- 442 14. Martin DP, Posada D, Crandall KA, Williamson C. A modified bootscan algorithm for  
443 automated identification of recombinant sequences and recombination breakpoints. *AIDS Res*  
444 *Hum Retroviruses*. 2005;21(1):98-102. doi: 10.1089/aid.2005.21.98. PubMed PMID: 15665649.
- 445 15. Smith JM. Analyzing the mosaic structure of genes. *Journal of molecular evolution*.  
446 1992;34(2):126-9. PubMed PMID: 1556748.
- 447 16. Posada D. Evaluation of methods for detecting recombination from DNA sequences:  
448 empirical data. *Molecular biology and evolution*. 2002;19(5):708-17. PubMed PMID: 11961104.
- 449 17. Gibbs MJ, Armstrong JS, Gibbs AJ. Sister-scanning: a Monte Carlo procedure for  
450 assessing signals in recombinant sequences. *Bioinformatics*. 2000;16(7):573-82. PubMed  
451 PMID: 11038328.
- 452 18. Huson DH, Bryant D. Application of phylogenetic networks in evolutionary studies.  
453 *Molecular biology and evolution*. 2006;23(2):254-67. Epub 2005/10/14. doi:  
454 10.1093/molbev/msj030. PubMed PMID: 16221896.
- 455 19. Zhan XY, Zhu QY. Molecular evolution of virulence genes and non-virulence genes in  
456 clinical, natural and artificial environmental *Legionella pneumophila* isolates. *PeerJ*.  
457 2017;5:e4114. Epub 2017/12/12. doi: 10.7717/peerj.4114. PubMed PMID: 29226035; PubMed  
458 Central PMCID: PMCPMC5719964.
- 459 20. Gao F, Chen C, Arab DA, Du Z, He Y, Ho SYWJE, et al. EasyCodeML: A visual tool for  
460 analysis of selection using CodeML. 2019.
- 461 21. Yang Z, Wong WS, Nielsen R. Bayes empirical bayes inference of amino acid sites under



- 462 positive selection. *Molecular biology and evolution*. 2005;22(4):1107-18. doi:  
463 10.1093/molbev/msi097. PubMed PMID: 15689528.
- 464 22. Kosakovsky Pond SL, Poon AFY, Velazquez R, Weaver S, Hepler NL, Murrell B, et al.  
465 HyPhy 2.5—A Customizable Platform for Evolutionary Hypothesis Testing Using Phylogenies.  
466 *Molecular biology and evolution*. 2019;37(1):295-9. doi: 10.1093/molbev/msz197 *Molecular*  
467 *Biology and Evolution*.
- 468 23. The PyMOL Molecular Graphics System, Version 1.5.X Schrödinger, LLC.
- 469 24. Narayanan K, Makino S. Cooperation of an RNA packaging signal and a viral envelope  
470 protein in coronavirus RNA packaging. *Journal of virology*. 2001;75(19):9059-67. Epub  
471 2001/09/05. doi: 10.1128/JVI.75.19.9059-9067.2001. PubMed PMID: 11533169; PubMed  
472 Central PMCID: PMCPMC114474.
- 473 25. Letko M, Marzi A, Munster V. Functional assessment of cell entry and receptor usage for  
474 SARS-CoV-2 and other lineage B betacoronaviruses. *Nat Microbiol*. 2020;5(4):562-9. Epub  
475 2020/02/26. doi: 10.1038/s41564-020-0688-y. PubMed PMID: 32094589; PubMed Central  
476 PMCID: PMCPMC7095430.
- 477 26. Voss D, Kern A, Traggiai E, Eickmann M, Stadler K, Lanzavecchia A, et al.  
478 Characterization of severe acute respiratory syndrome coronavirus membrane protein. *FEBS*  
479 *letters*. 2006;580(3):968-73. Epub 2006/01/31. doi: 10.1016/j.febslet.2006.01.026. PubMed  
480 PMID: 16442106; PubMed Central PMCID: PMCPMC7094741.
- 481 27. Tseng YT, Wang SM, Huang KJ, Lee AI, Chiang CC, Wang CT. Self-assembly of severe  
482 acute respiratory syndrome coronavirus membrane protein. *The Journal of biological chemistry*.  
483 2010;285(17):12862-72. Epub 2010/02/16. doi: 10.1074/jbc.M109.030270. PubMed PMID:

- 484 20154085; PubMed Central PMCID: PMCPMC2857088.
- 485 28. Wong MC, Javornik Cregeen SJ, Ajami NJ, Petrosino JF. Evidence of recombination in  
486 coronaviruses implicating pangolin origins of nCoV-2019. 2020:2020.02.07.939207. doi:  
487 10.1101/2020.02.07.939207 bioRxiv.
- 488 29. Wu Y. Strong evolutionary convergence of receptor-binding protein spike between COVID-  
489 19 and SARS-related coronaviruses. 2020:2020.03.04.975995. doi:  
490 10.1101/2020.03.04.975995. bioRxiv.
- 491 30. Wu A, Niu P, Wang L, Zhou H, Zhao X, Wang W, et al. Mutations, Recombination and  
492 Insertion in the Evolution of 2019-nCoV. 2020:2020.02.29.971101. doi:  
493 10.1101/2020.02.29.971101 bioRxiv.
- 494 31. Iceland patient infected by two strains. The Standard.  
495 2020;[https://www.thestandard.com.hk/section-news/section/11/217711/Iceland-patient--](https://www.thestandard.com.hk/section-news/section/11/217711/Iceland-patient--infected-by--two-strains)  
496 [infected-by--two-strains](https://www.thestandard.com.hk/section-news/section/11/217711/Iceland-patient--infected-by--two-strains).
- 497 32. Mallapaty S. How sewage could reveal true scale of coronavirus outbreak. Nature.  
498 2020;580(7802):176-7. Epub 2020/04/05. doi: 10.1038/d41586-020-00973-x. PubMed PMID:  
499 32246117.
- 500 33. Li X, Giorgi EE, Marichann MH, Foley B, Xiao C, Kong X-P, et al. Emergence of SARS-  
501 CoV-2 through Recombination and Strong Purifying Selection. 2020:2020.03.20.000885. doi:  
502 10.1101/2020.03.20.000885. bioRxiv.
- 503 34. Hughes AL, Hughes MA. More effective purifying selection on RNA viruses than in DNA  
504 viruses. Gene. 2007;404(1-2):117-25. Epub 2007/10/12. doi: 10.1016/j.gene.2007.09.013.  
505 PubMed PMID: 17928171; PubMed Central PMCID: PMCPMC2756238.

- 506 35. Pond SLK, Muse SV. HyPhy: Hypothesis Testing Using Phylogenies: Springer New York;  
507 2005. 676-9 p.
- 508 36. Pond SL, Scheffler K, Gravenor MB, Poon AF, Frost SD. Evolutionary fingerprinting of  
509 genes. *Molecular biology and evolution*. 2010;27(3):520-36. Epub 2009/10/30. doi:  
510 10.1093/molbev/msp260. PubMed PMID: 19864470; PubMed Central PMCID:  
511 PMCPMC2877558.
- 512 37. Kosakovsky Pond SL, Frost SD. Not so different after all: a comparison of methods for  
513 detecting amino acid sites under selection. *Molecular biology and evolution*. 2005;22(5):1208-  
514 22. Epub 2005/02/11. doi: 10.1093/molbev/msi105. PubMed PMID: 15703242.
- 515 38. Wan Y, Shang J, Graham R, Baric RS, Li F. Receptor Recognition by the Novel  
516 Coronavirus from Wuhan: an Analysis Based on Decade-Long Structural Studies of SARS  
517 Coronavirus. *Journal of virology*. 2020;94(7). Epub 2020/01/31. doi: 10.1128/JVI.00127-20.  
518 PubMed PMID: 31996437; PubMed Central PMCID: PMCPMC7081895.
- 519 39. Walls AC, Park YJ, Tortorici MA, Wall A, McGuire AT, Velesler D. Structure, Function, and  
520 Antigenicity of the SARS-CoV-2 Spike Glycoprotein. *Cell*. 2020;181(2):281-92 e6. Epub  
521 2020/03/11. doi: 10.1016/j.cell.2020.02.058. PubMed PMID: 32155444; PubMed Central  
522 PMCID: PMCPMC7102599.
- 523 40. Clement M, Snell Q, Walker P, Posada D, Crandall KJP, Distributed Processing  
524 Symposium IP. *TCS: Estimating gene genealogies*. 2002;2:184.
- 525 41. Bubac CM, Spellman GMJTAOA. How connectivity shapes genetic structure during range  
526 expansion: Insights from the Virginia's Warbler. 2016;(2):2.
- 527 42. Duxbury EM, Day JP, Maria Vespasiani D, Thuringer Y, Tolosana I, Smith SC, et al. Host-

- 528 pathogen coevolution increases genetic variation in susceptibility to infection. *eLife*. 2019;8.  
529 Epub 2019/05/01. doi: 10.7554/eLife.46440. PubMed PMID: 31038124; PubMed Central  
530 PMCID: PMC6491035.
- 531 43. Meyerson NR, Sawyer SL. Two-stepping through time: mammals and viruses. *Trends in*  
532 *microbiology*. 2011;19(6):286-94. Epub 2011/05/03. doi: 10.1016/j.tim.2011.03.006. PubMed  
533 PMID: 21531564; PubMed Central PMCID: PMC3567447.
- 534 44. Lu G, Wang Q, Gao GF. Bat-to-human: spike features determining 'host jump' of  
535 coronaviruses SARS-CoV, MERS-CoV, and beyond. *Trends in microbiology*. 2015;23(8):468-  
536 78. Epub 2015/07/25. doi: 10.1016/j.tim.2015.06.003. PubMed PMID: 26206723; PubMed  
537 Central PMCID: PMC7125587.
- 538 45. Bestle D, Heindl MR, Limburg H, van TVL, Pilgram O, Moulton H, et al. TMPRSS2 and  
539 furin are both essential for proteolytic activation and spread of SARS-CoV-2 in human airway  
540 epithelial cells and provide promising drug targets. 2020:2020.04.15.042085. doi:  
541 10.1101/2020.04.15.042085. bioRxiv.
- 542 46. Ou X, Liu Y, Lei X, Li P, Mi D, Ren L, et al. Characterization of spike glycoprotein of SARS-  
543 CoV-2 on virus entry and its immune cross-reactivity with SARS-CoV. *Nature communications*.  
544 2020;11(1):1620. Epub 2020/03/30. doi: 10.1038/s41467-020-15562-9. PubMed PMID:  
545 32221306; PubMed Central PMCID: PMC7100515.
- 546 47. Gui M, Song W, Zhou H, Xu J, Chen S, Xiang Y, et al. Cryo-electron microscopy structures  
547 of the SARS-CoV spike glycoprotein reveal a prerequisite conformational state for receptor  
548 binding. *Cell research*. 2017;27(1):119-29. Epub 2016/12/23. doi: 10.1038/cr.2016.152.  
549 PubMed PMID: 28008928; PubMed Central PMCID: PMC5223232.

550 48. Wrapp D, Wang N, Corbett KS, Goldsmith JA, Hsieh CL, Abiona O, et al. Cryo-EM  
551 structure of the 2019-nCoV spike in the prefusion conformation. *Science*.  
552 2020;367(6483):1260-3. Epub 2020/02/23. doi: 10.1126/science.abb2507. PubMed PMID:  
553 32075877.

554 49. Walls AC, Xiong X, Park YJ, Tortorici MA, Snijder J, Quispe J, et al. Unexpected Receptor  
555 Functional Mimicry Elucidates Activation of Coronavirus Fusion. *Cell*. 2019;176(5):1026-39 e15.  
556 Epub 2019/02/05. doi: 10.1016/j.cell.2018.12.028. PubMed PMID: 30712865; PubMed Central  
557 PMCID: PMC6751136.

558 50. Walls AC, Tortorici MA, Snijder J, Xiong X, Bosch BJ, Rey FA, et al. Tectonic  
559 conformational changes of a coronavirus spike glycoprotein promote membrane fusion.  
560 *Proceedings of the National Academy of Sciences of the United States of America*.  
561 2017;114(42):11157-62. Epub 2017/10/27. doi: 10.1073/pnas.1708727114. PubMed PMID:  
562 29073020; PubMed Central PMCID: PMC5651768.

563 51. Xiong X, Qu K, Ciazynska KA, Hosmillo M, Carter AP, Ebrahimi S, et al. A thermostable,  
564 closed, SARS-CoV-2 spike protein trimer. 2020:2020.06.15.152835. doi:  
565 10.1101/2020.06.15.152835. bioRxiv.  
566  
567

## 568 **Supporting information**

569 **S1 Table.** SARS-CoV-2 isolates information.

570 **S2 Table.** Log-likelihood values and parameter estimates for the SARS-CoV-2 N gene sequences.

571 **S3 Table.** Detailed information of SARS-CoV-2 isolates with full length sequence of S gene. The  
572 data are organized by weekly.

573 **S4 Table.** Detailed information of SARS-CoV-2 isolates with full length sequence of S gene. The  
574 data are organized by panels. Each panel contains 200-300 isolates by combining isolates from  
575 several days.

576 **S1 Fig. The evolutionary relationship of N alleles. A.** Phylogenetic tree of N gene based on  
577 nucleotide sequences of 131 alleles. The evolutionary history is inferred using the Maximum  
578 Likelihood method and Tamura-Nei model. The tree is drawn to scale, with branch lengths measured  
579 in the number of substitutions per site. Bootstrap values more than 0.5 are shown. **B.** Parsimony  
580 network of SARS-CoV-2 N gene haplotype (allele) diversity obtained from 3090 isolates worldwide.  
581 Each oblique line linking between haplotypes (haplotype name is shown as its representative isolate  
582 name) represents one mutational difference. The ancestral haplotype, or root of the network, is  
583 labeled with a square, and represent haplotype name is marked red.

584

585 **S2 Fig. Evolutionary relationship of S alleles with or without L5F mutation. A.** Phylogenetic  
586 tree of S gene based on nucleotide sequences of 173 alleles. Each clade is highlighted with different  
587 color. Alleles are shown with their representative isolate names, and alleles with L5F mutation are  
588 highlighted in blue. Bootstrap values more than 0.5 are shown. **B.** Parsimony network of SARS-  
589 CoV-2 S gene haplotype (allele) diversity obtained from 3090 isolates worldwide. Each oblique line

590 linking between haplotypes (haplotype name is shown as its representative isolate name) represents  
591 one mutational difference. Unlabeled nodes (Gray circle) indicate inferred steps have not found in  
592 the sampled populations yet. The ancestral haplotype, or root of the network, is labeled with a square,  
593 and represent haplotype name is marked green or red. The blue nodes indicate haplotypes with L5F  
594 mutation. Dotted boxes indicate major haplotype groups. Haplotypes include in red dotted boxes  
595 are with D614G mutation while those included in black dotted boxes are without D614G mutation.

596

597 **S3 Fig. Positive selection analysis of S gene codon 5 by Evolutionary Fingerprinting method.**

598 Log (Bayes Factor) for positive selection at codon 5 of S gene and its frequencies. The cut-off value  
599 for the Bayes factor (BF) in the Evolutionary Fingerprinting method was set at 25 to reflect a positive  
600 selection at a given site (Posterior probability>0.95).  $\Pr \{BF>25\}$  indicates posterior probability of  
601 Bayes Factor >25.

602

603

## 604 **Figure legends**

605

606 **Figure 1.** Phylogenetic tree of E (A), M (B), N (C), and S(D) proteins of SARS-CoV-2. Major  
607 clades are highlighted with different color. The tree shows topology of the protein of each allele,  
608 named by their representative isolates.

609

610 **Figure 2.** Reticulate network trees of E (A), M (B), N (C) and S (D) alleles of SARS-CoV-2  
611 analyzed by the neighbor-net algorithm of SplitsTree4. Scale bars indicate number of substitutions  
612 per site. All internal nodes represent hypothetical ancestral alleles and edges that correspond to  
613 reticulate events such as recombination. Red arrows indicate edges. Because there are too few  
614 informative characters to use the Phi test for E and M genes, *p-values* of Phi test of N and S genes  
615 are shown.

616

617 **Figure 3.** Tajima's D, Fu and Li's D\* and F\* test for the four structural gene alleles of SARS-  
618 CoV-2. \* $p < 0.05$ ; \*\* $p < 0.01$ ; \*\*\* $p < 0.001$

619

620

621 **Figure 4. Positive selection analysis of S gene codons by IFEL and Evolutionary**

622 **Fingerprinting methods.** **A.** Diagram of selection analysis result of S codons by IFEL method.

623 Asterisk indicates the positive selection site with statistical significance ( $p < 0.01$ ). **B.** Log (Bayes

624 Factor) for positive selection at codon 614 of S gene and its frequencies. The cut-off value for the

625 Bayes factor (BF) in the Evolutionary Fingerprinting method was set at 25 to reflect a positive



626 selection at a given site (Posterior probability>0.95). Pr {BF>25} indicates posterior probability of  
627 Bayes Factor >25.

628

629 **Figure 5. Evolutionary relationship of S alleles with or without D614G mutation. A.**

630 Phylogenetic tree of S gene based on nucleotide sequences of 173 alleles. The evolutionary history  
631 is inferred using the Maximum Likelihood method and Tamura-Nei model. The tree is drawn to  
632 scale, with branch lengths measured in the number of substitutions per site. Each clade is highlighted  
633 with different color. Alleles are shown with their representative isolate names, and alleles with  
634 D614G mutation are highlighted in red. Bootstrap values more than 0.5 are shown. **B.** Parsimony  
635 network of SARS-CoV-2 S gene haplotype (allele) diversity obtained from 3090 isolates worldwide.  
636 Each oblique line linking between haplotypes (haplotype name is shown as its representative isolate  
637 name) represents one mutational difference. Unlabeled nodes (Gray circle) indicate inferred steps  
638 have not found in the sampled populations yet. The ancestral haplotype, or root of the network, is  
639 labeled with a square, and represent haplotype name is marked green or red. The red nodes indicate  
640 haplotypes with D614G mutation, while green or black nodes indicate haplotypes without D614G  
641 mutation. Dotted boxes indicate major haplotype groups.

642

643 **Figure 6. Expansion of S alleles with D614G mutation during SARS-CoV-2 human to human**

644 transmission. **A.** Percentage of SARS-CoV-2 isolates carrying the alleles of D614G mutation in  
645 each week collected. **B.** Frequencies of D614G mutation in the S gene in each period of time (Four  
646 to five weeks' data are combined). **C.** Percentage of SARS-CoV-2 isolates carrying the alleles with

647 D614G mutation in each period of time. **D.** Frequencies of D614G mutation in the S gene in each  
648 period of time. \* $p < 0.05$ ; \*\* $p < 0.01$ .

649

650

651 **Figure 7.** The structure of the S protein of SARS-CoV-2 and potential influence of D614G mutation  
652 on its structural change. **A.** Schematic of the primary structure of SARS-CoV-2 S protein colored  
653 by domains. Some boundary-residues are listed. The S1/S2 cleavage sites are indicated by arrows.

654 RBD: receptor binding domain; RBM: receptor of binding motif; FP: fusion peptide, HR1/2: heptad  
655 repeat 1/2; TM: transmembrane domain; CT: cytoplasmic tail; NTD: N-terminal domain; CTD: C-

656 terminal domain; SD1: subdomain 1; SD2: subdomain 2. The structure of the S protein trimer of

657 SARS-CoV-2 and potential influence of D614G mutation on its structural change. **B.**

658 Experimentally determined structure of SARS-CoV-2 S protein trimer (PDB ID is 6VSB and the  
659 amino acid sequences is the same as WH01 isolate). **C.** D614-K854 inter-monomer salt bridge. **D.**

660 G614-K854 inter-monomer salt bridge. The distance of the salt bridge is increased from 2.6 to 5.2

661 Å in D614G mutation as shown.

662

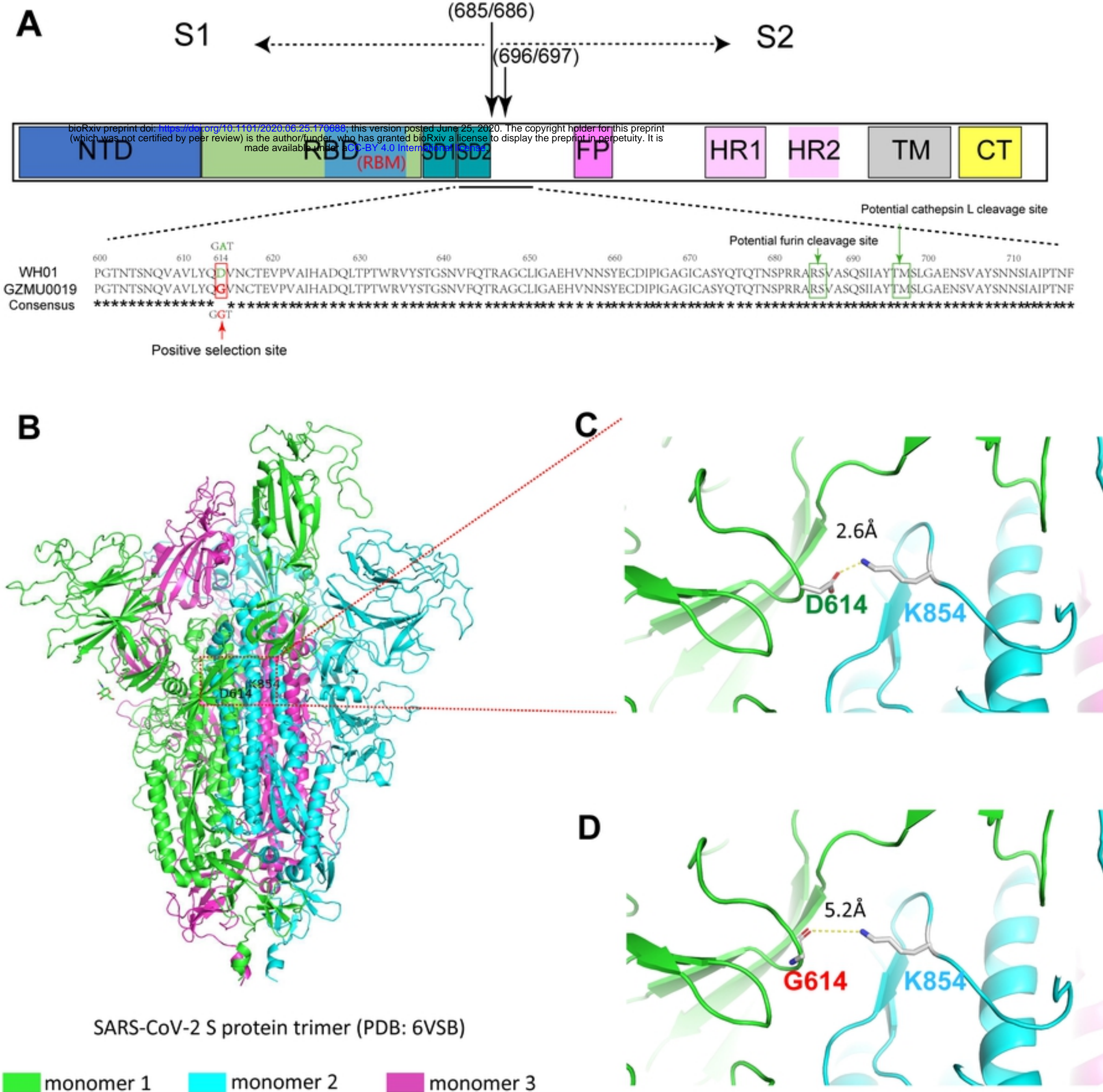


Figure 7



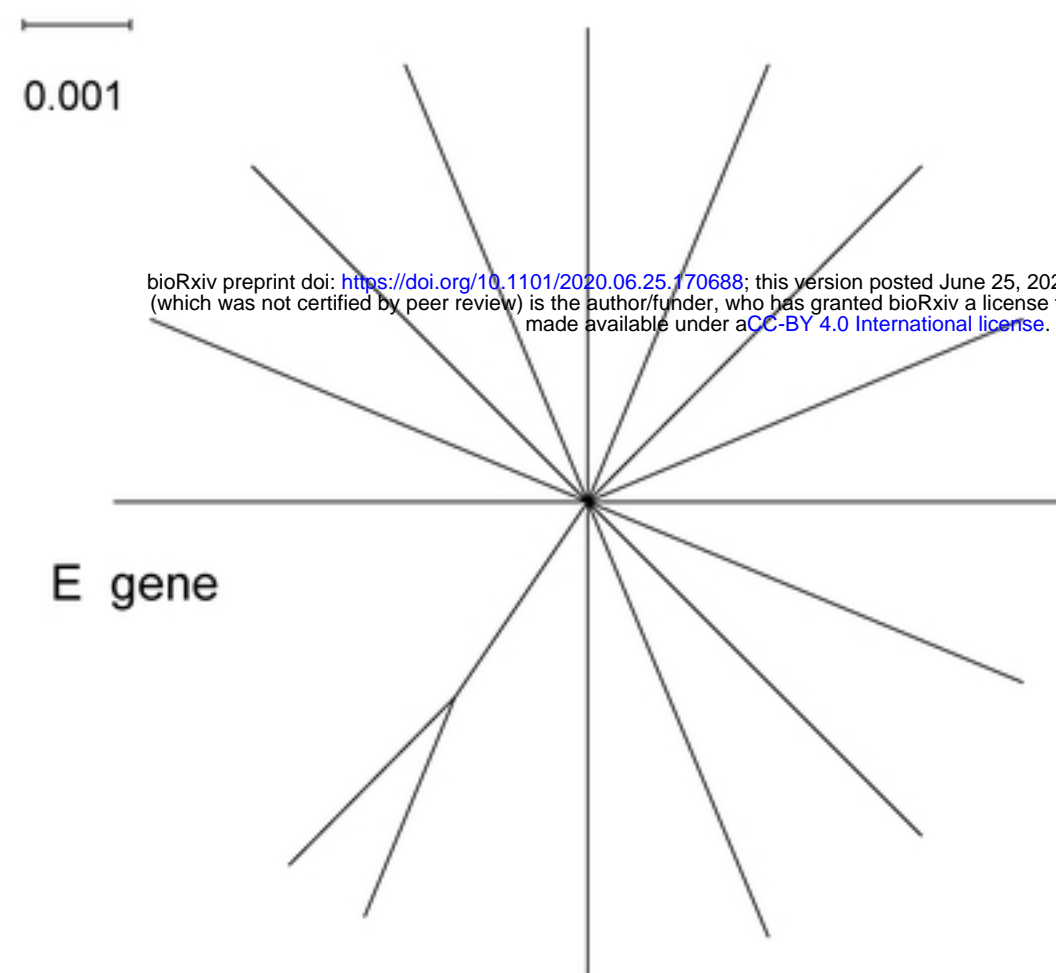
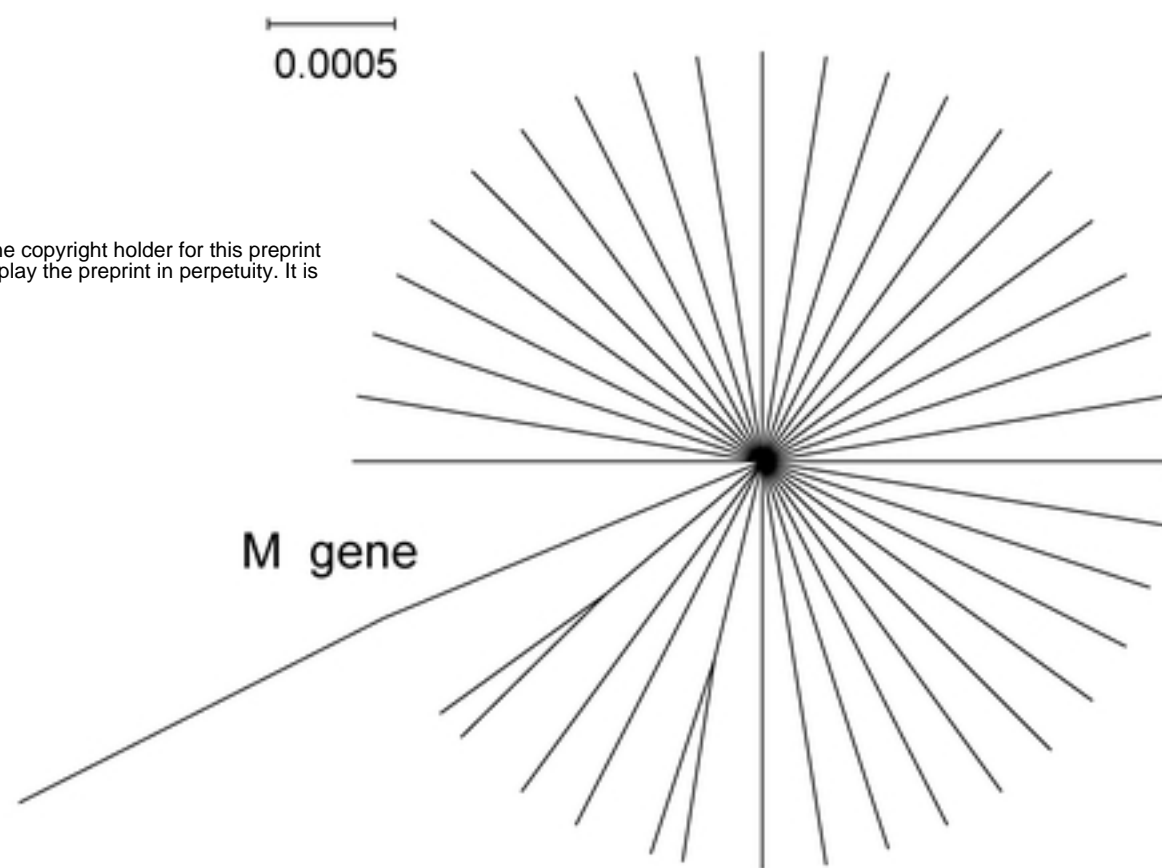
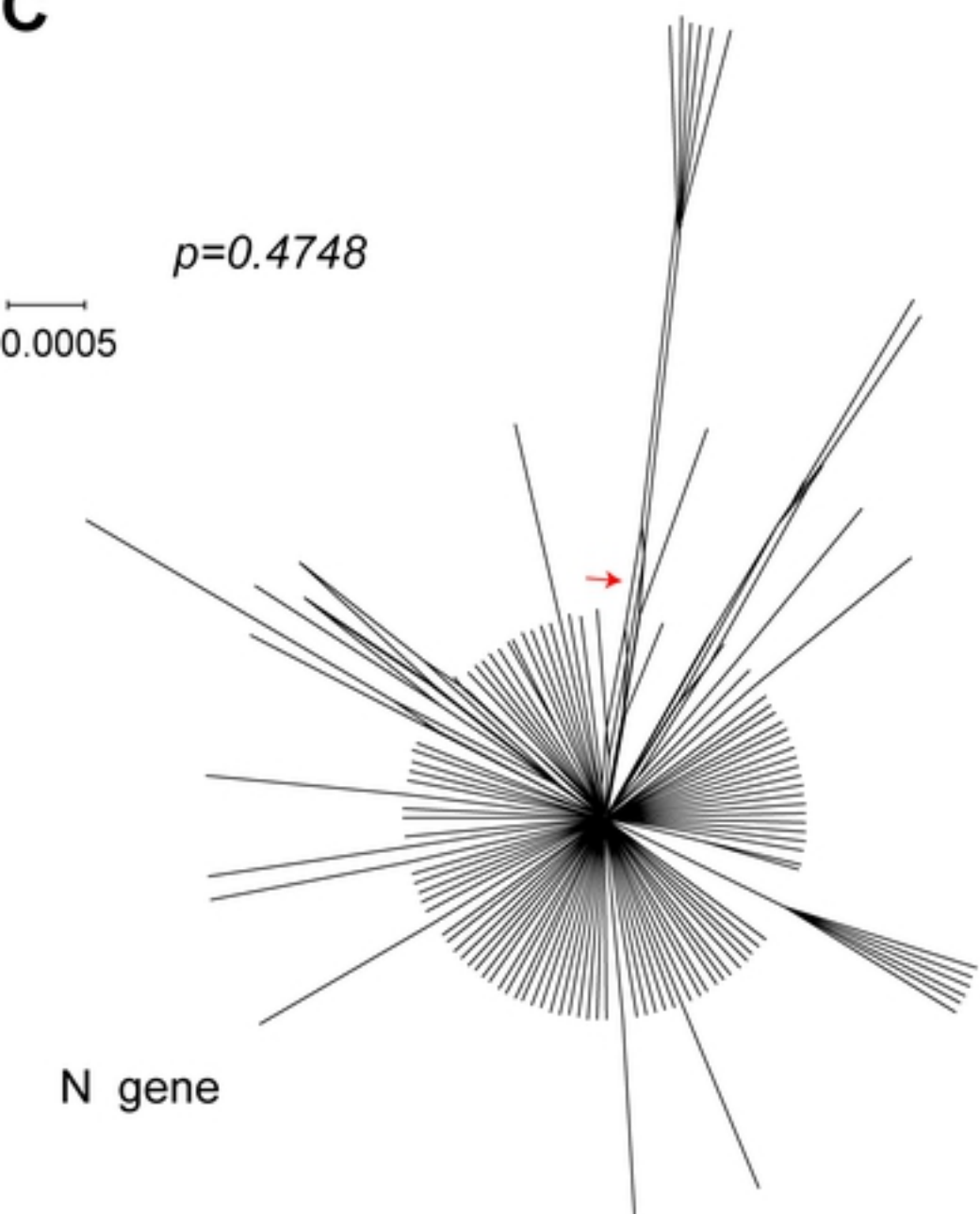
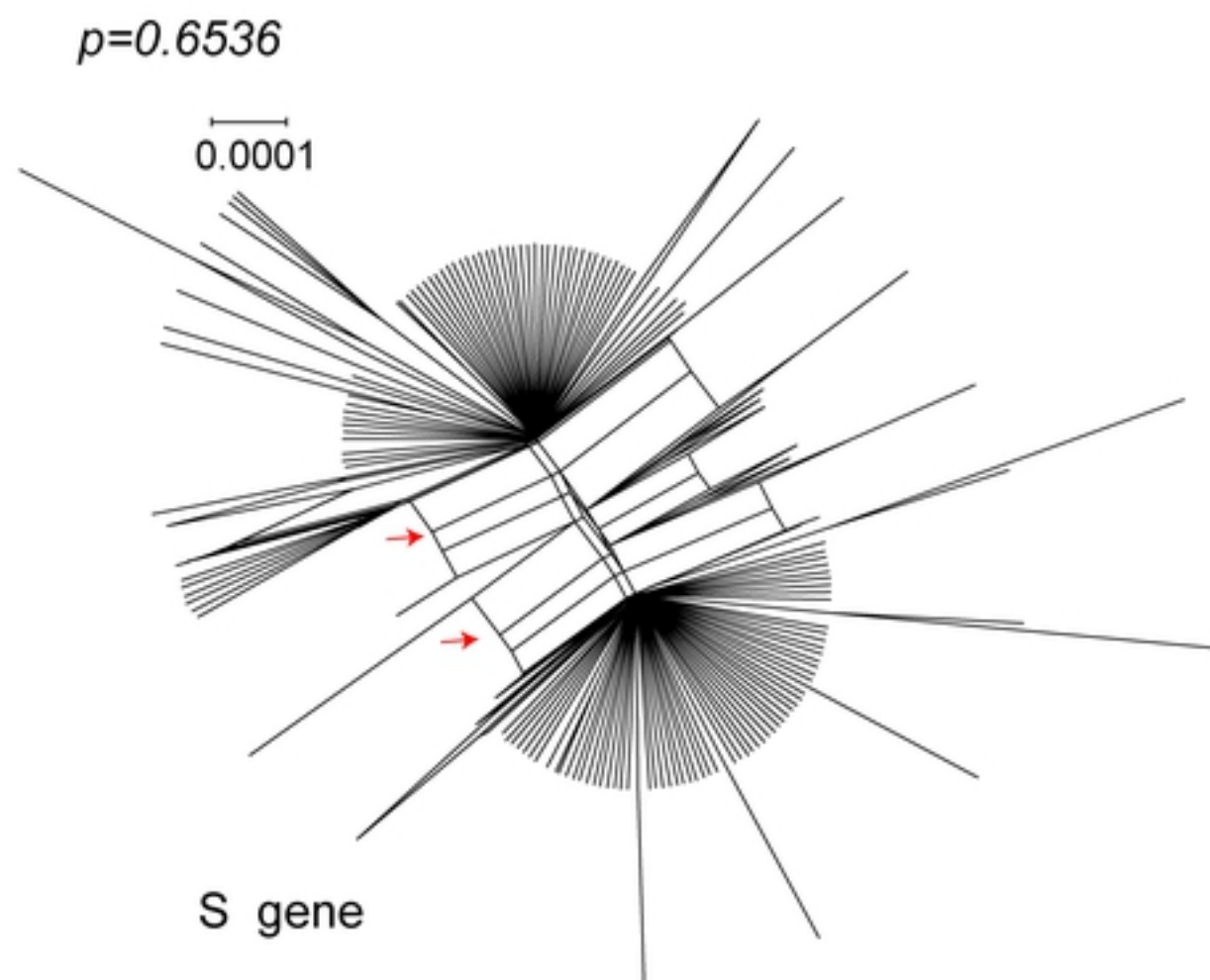
**A****B****C****D**

Figure 2

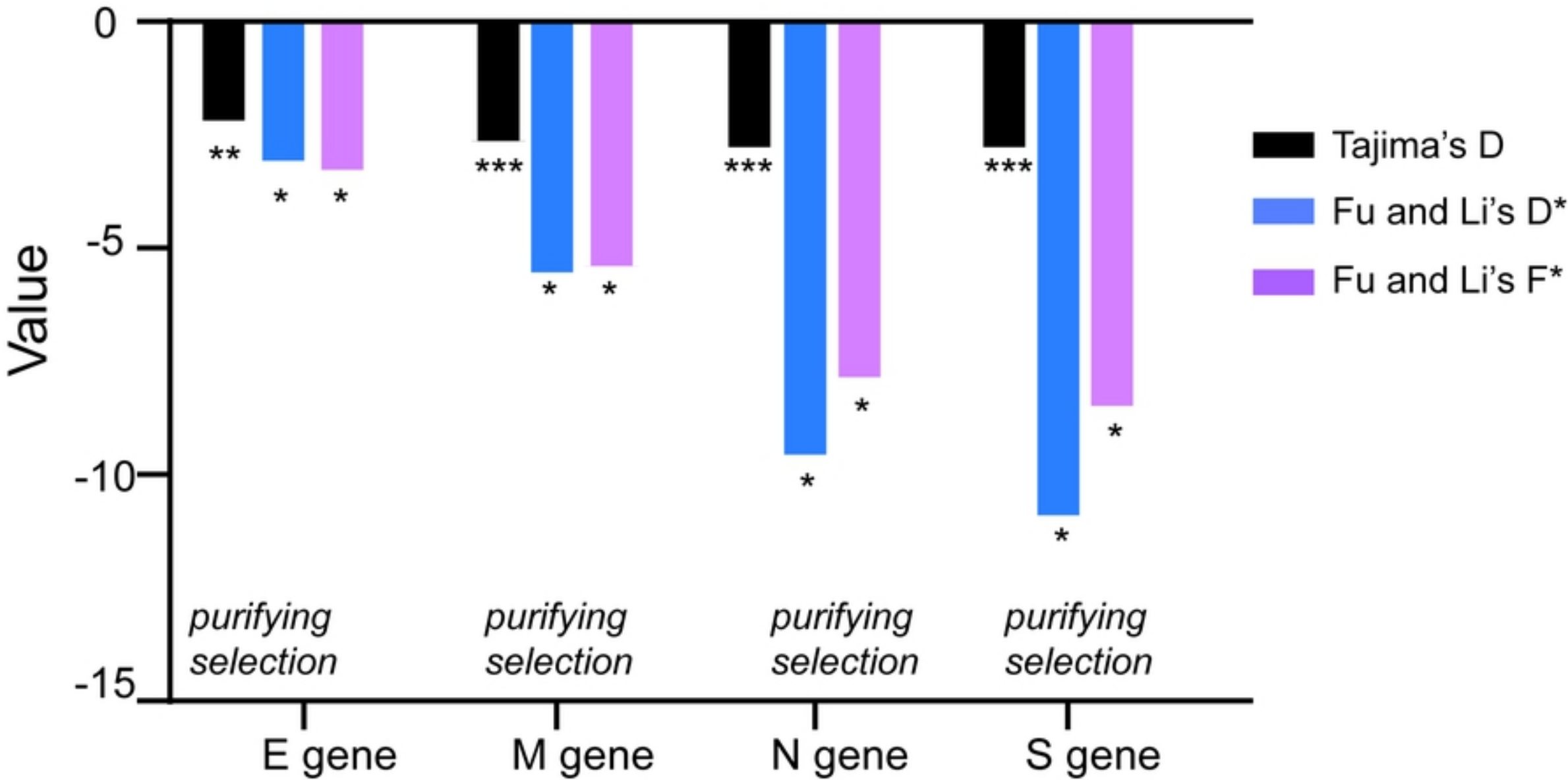


Figure 3

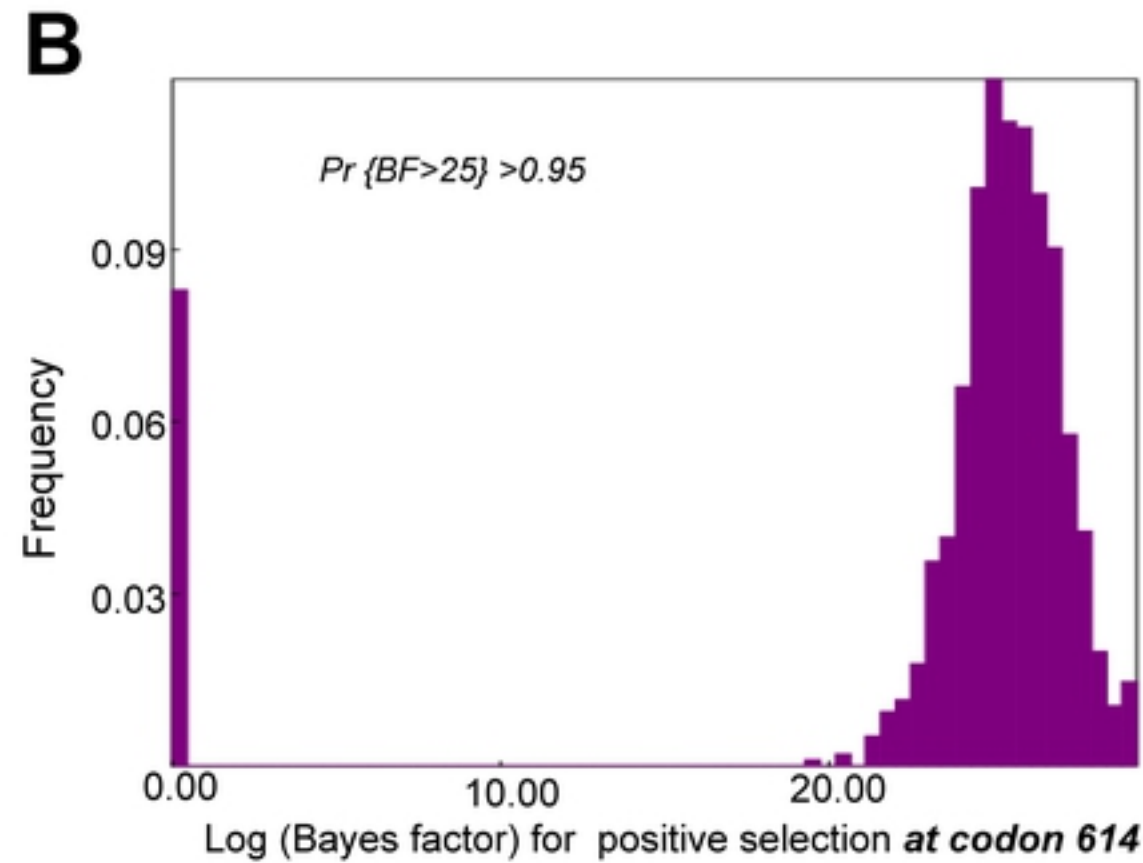
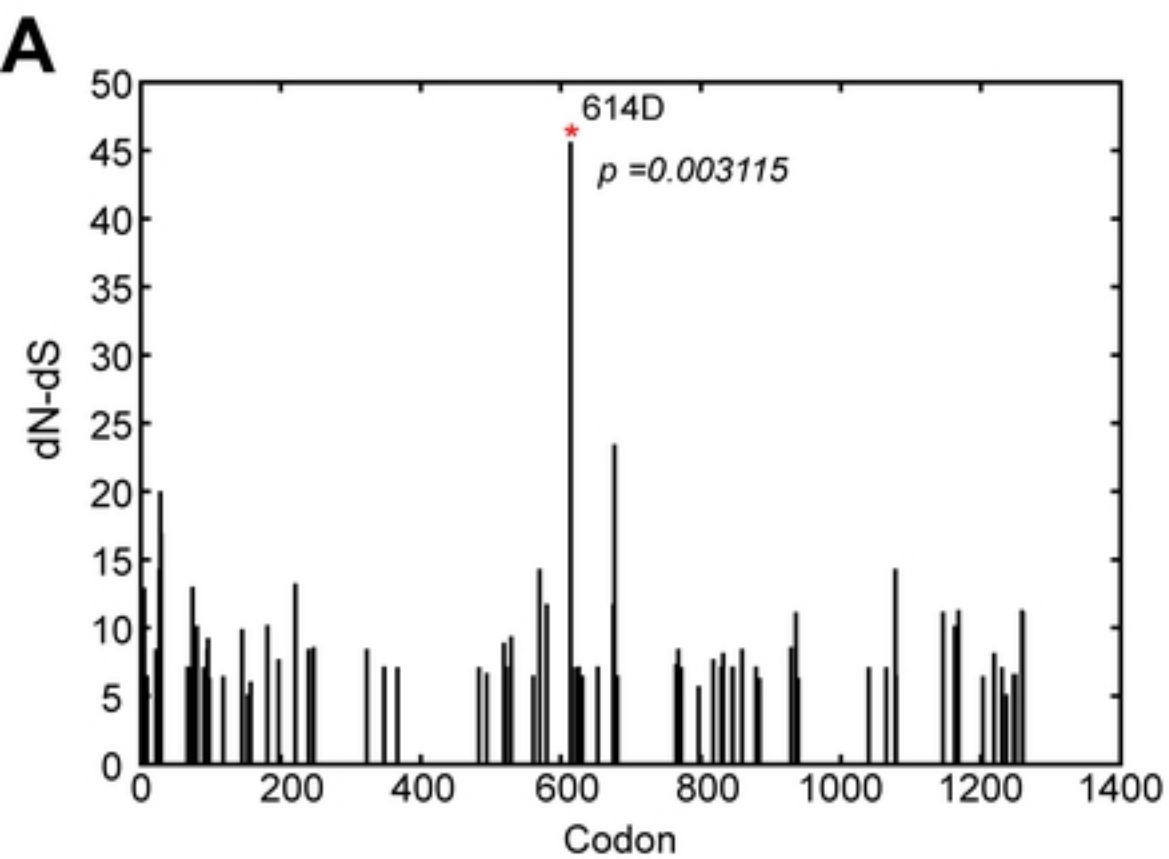


Figure 4

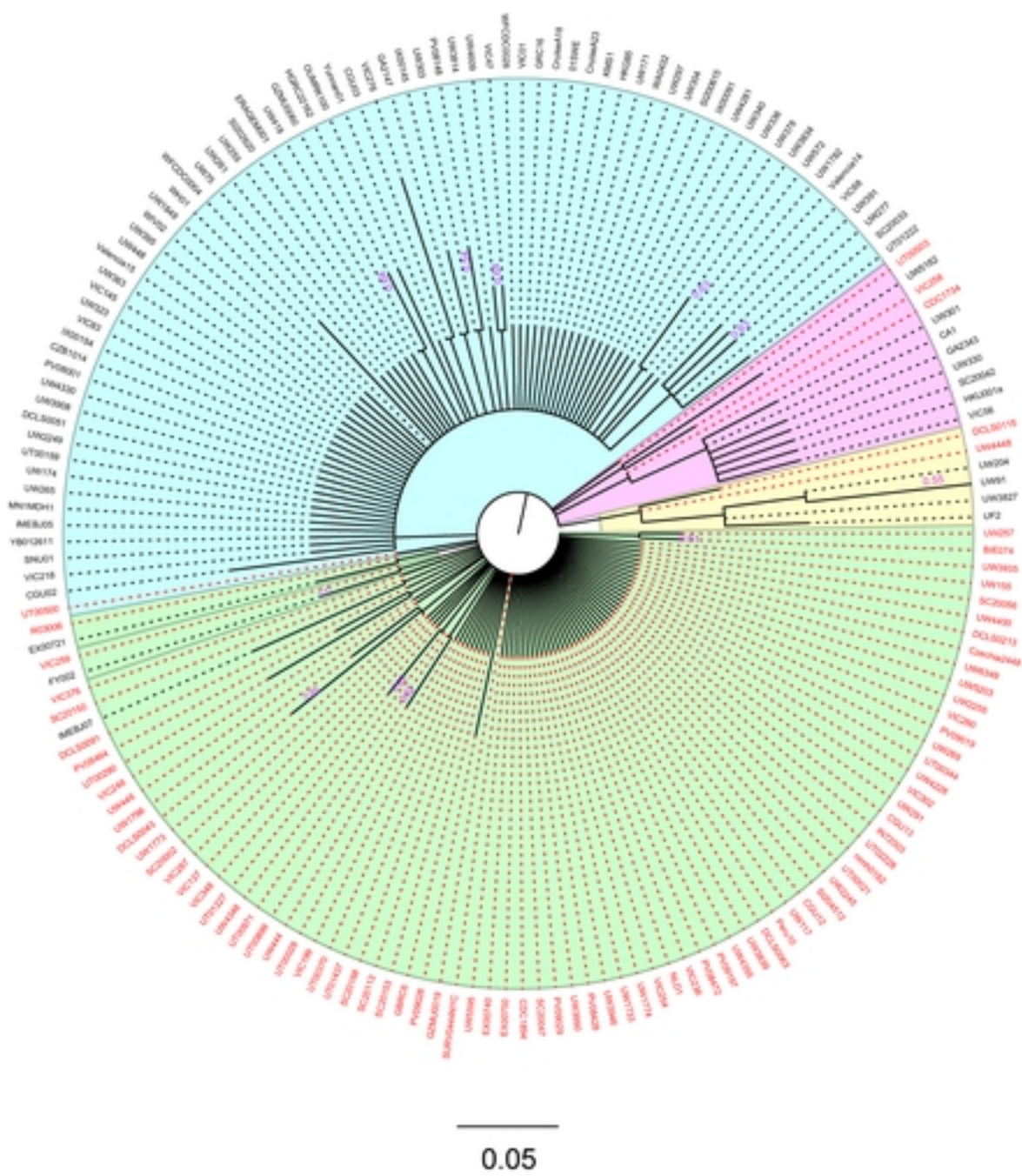
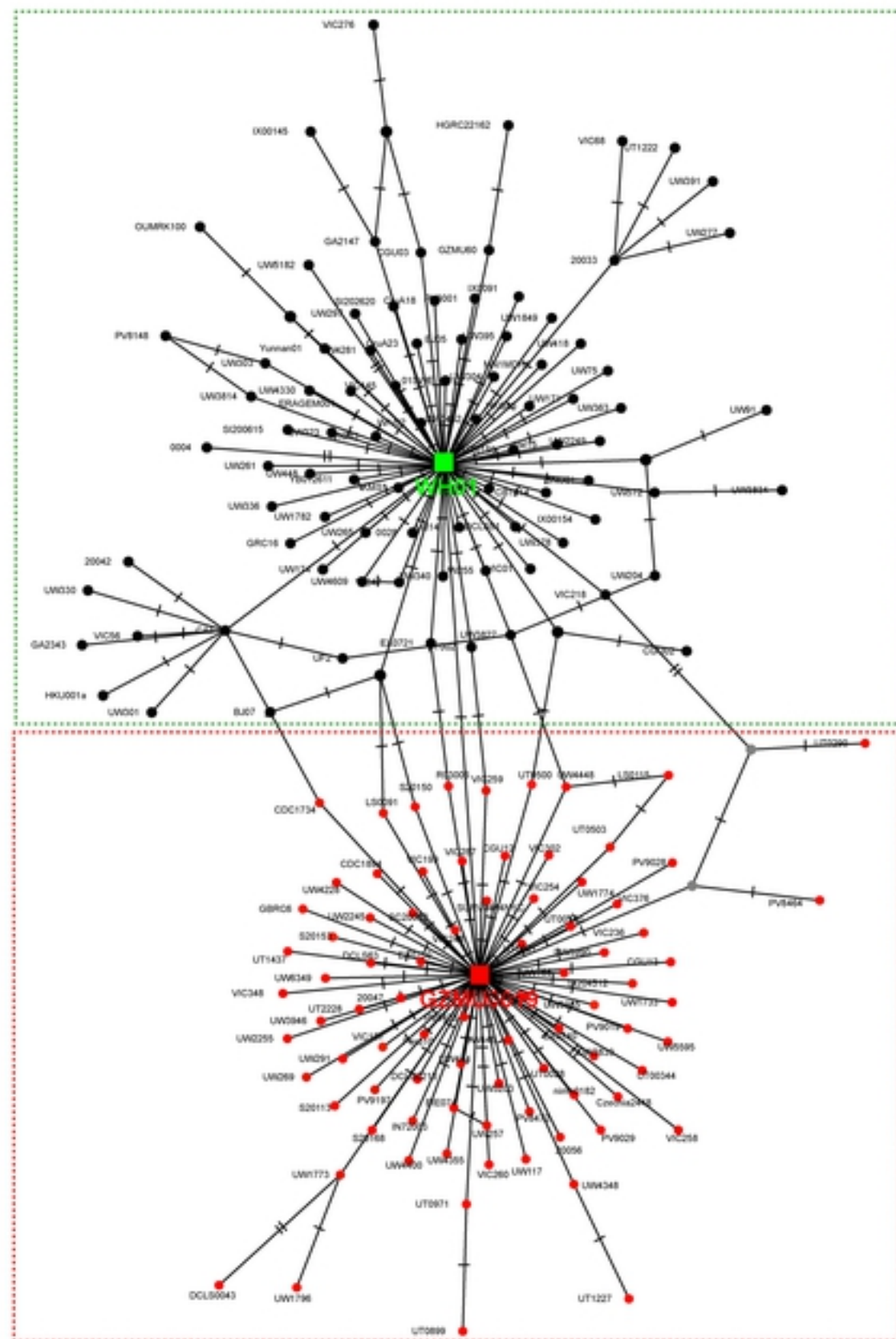
**A****B**

Figure 5



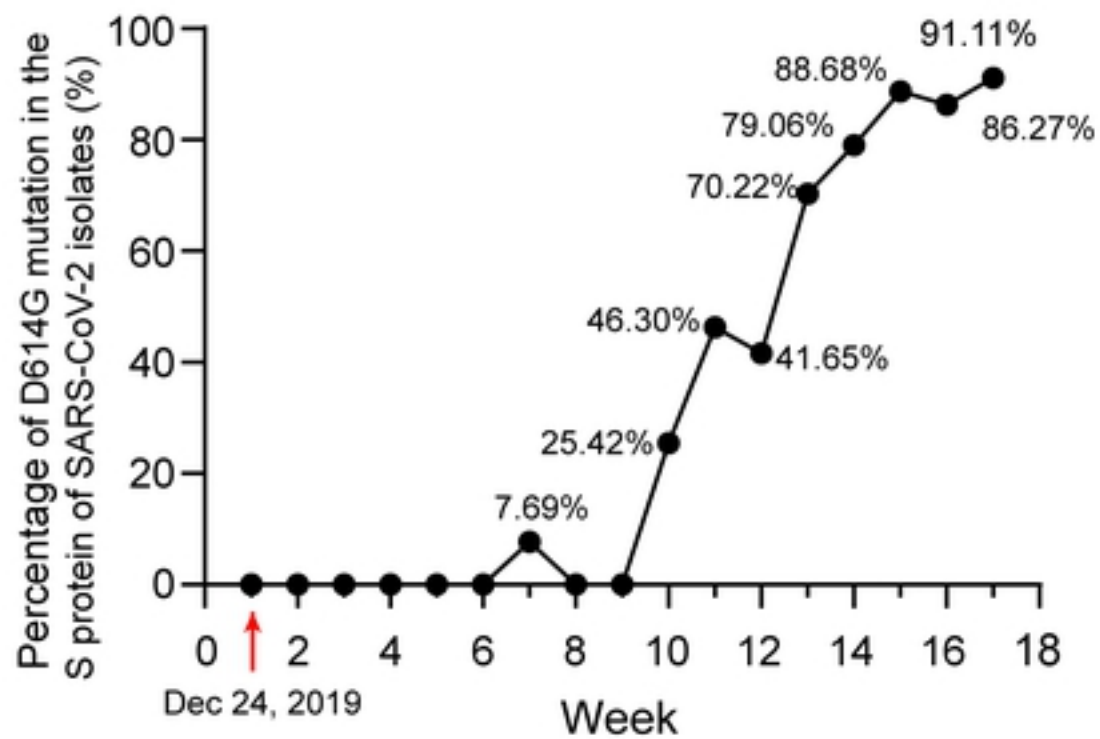
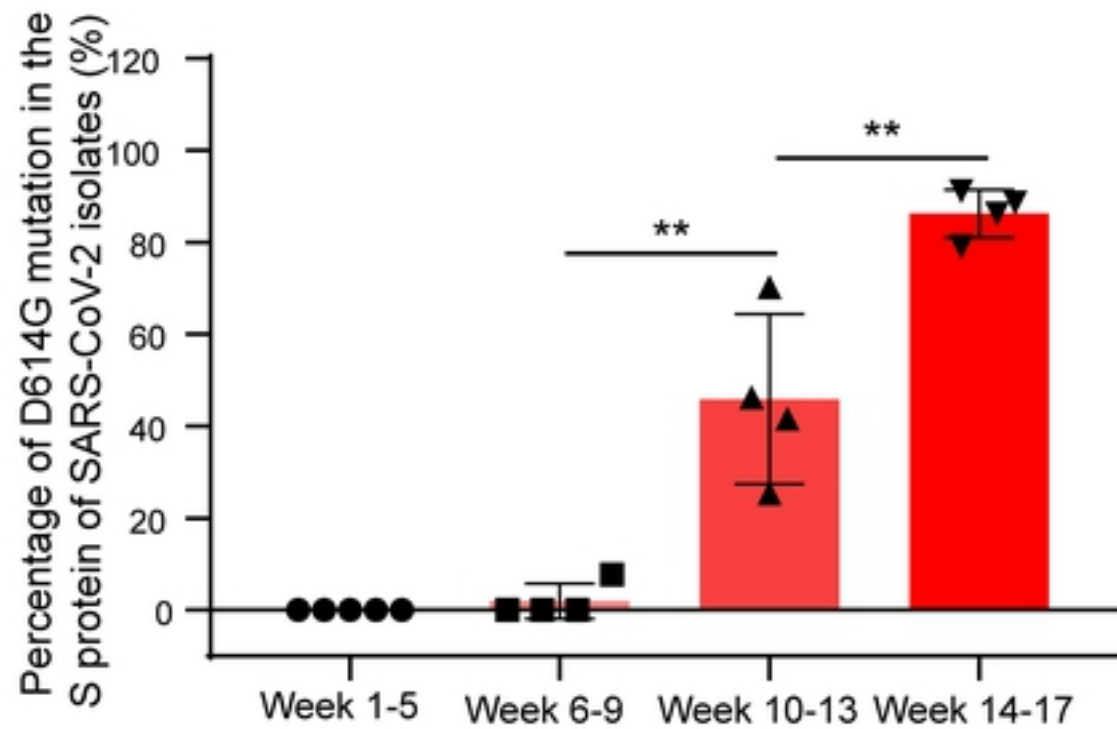
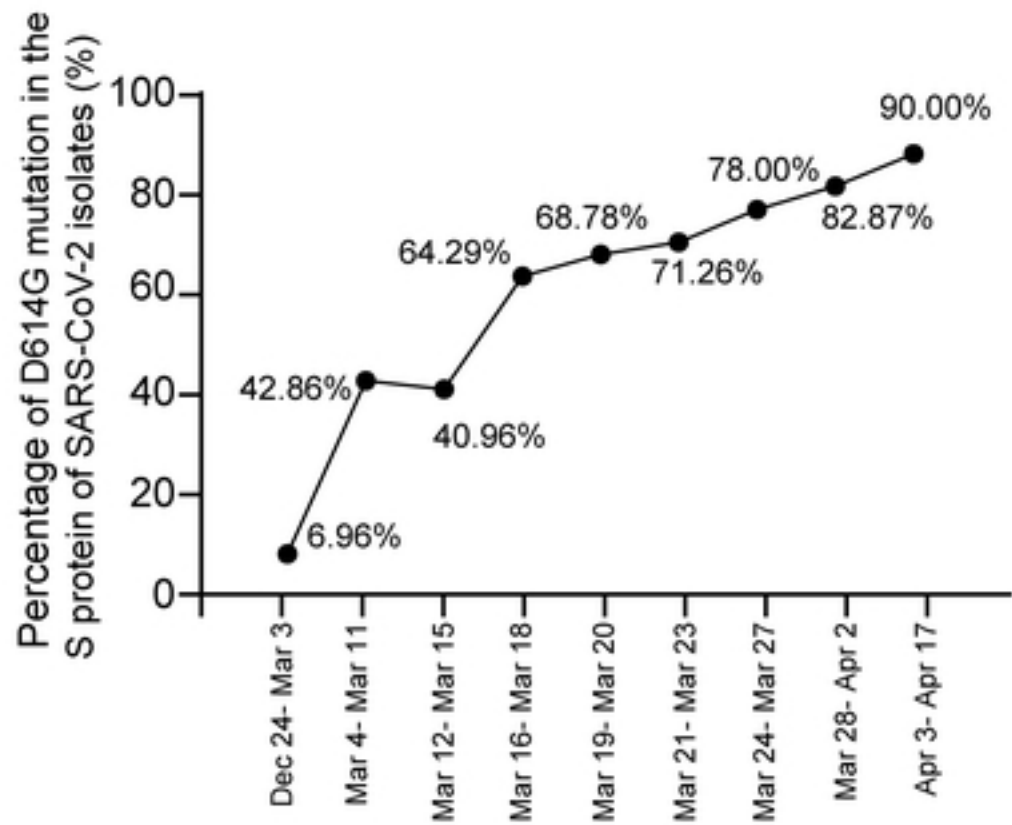
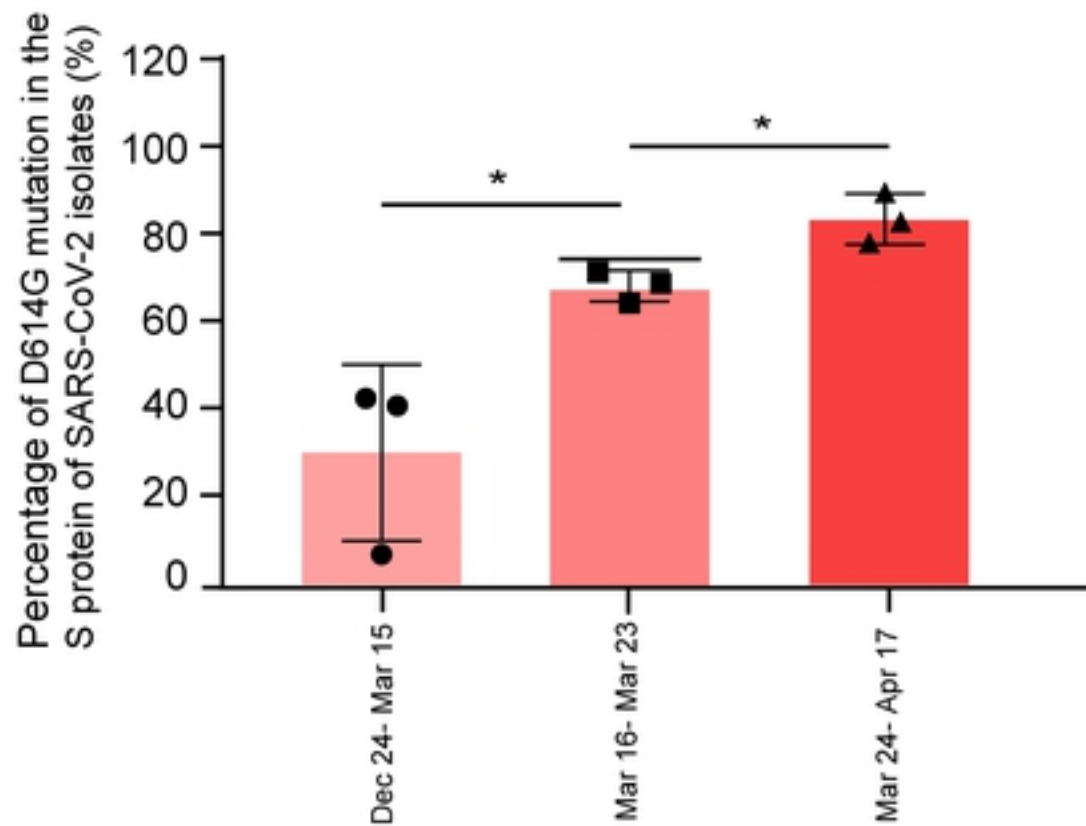
**A****B****C****D**

Figure 6