

1 **Multi-site co-mutations and 5'UTR CpG immunity escape drive** 2 **the evolution of SARS-CoV-2**

3 Jingsong Zhang¹, Junyan Kang^{2,3}, Mofang Liu^{2,3}, Benhao Han¹, Li Li⁴, Yongqun He⁵, Zhigang Yi^{6,7*},
4 Luonan Chen^{1,8,9,10*}

5
6 ¹Key Laboratory of Systems Biology, State Key Laboratory of Cell Biology, Shanghai Institute of
7 Biochemistry and Cell Biology, Center for Excellence in Molecular Cell Science, Chinese Academy
8 of Sciences, Shanghai 200031, China

9 ²State Key Laboratory of Molecular Biology, Shanghai Key Laboratory of Molecular Andrology,
10 Shanghai Institute of Biochemistry and Cell Biology, Center for Excellence in Molecular Cell
11 Science, Chinese Academy of Sciences, Shanghai 200031, China

12 ³University of Chinese Academy of Sciences, Shanghai 200031, China

13 ⁴Department of Genetics, Harvard Medical School, Boston, MA 02115, USA

14 ⁵University of Michigan Medical School, Ann Arbor, MI 48109, USA

15 ⁶Shanghai Public Health Clinical Center, Fudan University, Shanghai 201508, China

16 ⁷Key Laboratory of Medical Molecular Virology (MOE/NHC/CAMS), School of Basic Medical
17 Sciences, Shanghai Medical College, Fudan University, Shanghai 200032, China

18 ⁸School of Life Science and Technology, ShanghaiTech University, Shanghai 201210, China.

19 ⁹Key Laboratory of Systems Biology, Hangzhou Institute for Advanced Study, University of Chinese
20 Academy of Sciences, Hangzhou 310024, China

21 ¹⁰Center for Excellence in Animal Evolution and Genetics, Chinese Academy of Sciences, Kunming
22 650223, China.

23 *Corresponding author: Email: lnchen@sibs.ac.cn(L.N.C.); zgyi@fudan.edu.cn(Z.G.Y.)

29 **ABSTRACT**

30 **The SARS-CoV-2 infected cases and the caused mortalities have been surging since the**
31 **COVID-19 pandemic. Viral mutations emerge during the virus circulating in the population,**
32 **which is shaping the viral infectivity and pathogenicity. Here we extensively analyzed 6698**
33 **SARS-CoV-2 whole genome sequences with specific sample collection dates in NCBI database.**
34 **We found that four mutations, i.e., 5'UTR_c-241-t, NSP3_c-3037-t, NSP12_c-14408-t, and**
35 **S_a-23403-g, became the dominant variants and each of them represented nearly 100% of all**
36 **virus sequences since the middle May, 2020. Notably, we found that co-occurrence rates of**
37 **three significant multi-site co-mutational patterns, i.e., (i) S_a-23403-g, NSP12_c-14408-t,**
38 **5'UTR_c-241-t, NSP3_c-3037-t, and ORF3a_c-25563-t; (ii) ORF8_t-28144-c, NSP4_c-8782-t,**
39 **NSP14_c-18060-t, NSP13_a-17858-g, and NSP13_c-17747-t; and (iii) N_g-28881-a, N_g-28882-a,**
40 **and N_g-28883-c, reached 66%, 90%, and nearly 100% of recent sequences, respectively.**
41 **Moreover, we found significant decrease of CpG dinucleotide at positions 241(c)-242(g) in the**
42 **5'UTR during the evolution, which was verified as a potential target of human zinc finger**
43 **antiviral protein (ZAP). The four dominant mutations, three significant multi-site co-mutations,**
44 **and the potential escape mutation of ZAP-target in 5'UTR region contribute to the rapid**
45 **evolution of SARS-CoV-2 virus in the population, thus shaping the viral infectivity and**
46 **pathogenicity. This study provides valuable clues and frameworks to dissect the viral**
47 **replication and virus-host interactions for designing effective therapeutics.**

48

49 **INTRODUCTION**

50 Since the outbreak of COVID-19 in December 2019, it has been pandemic in over 200 countries.
51 The infected cases and the mortalities have been surging, which is an ongoing threat to the public
52 health (1, 2). COVID-19 is caused by infection with a novel coronavirus SARS-CoV-2 (3-5). Even
53 though as a coronavirus, SARS-CoV-2 has genetic proofreading mechanisms (6-8), the persistent
54 natural selective pressure in the population drives the virus to gradually accumulate favorable
55 mutations (6, 9). Considerable attention is given to the mutation and evolution of SARS-CoV-2, for
56 that viral mutations have important impact on the infection and pathogenicity of viruses (10). The
57 beneficial mutants can better evolve and adapt to host (9), either strengthening or weakening the
58 infectivity and pathogenicity. In addition, the variants may generate drug resistance and shrink the
59 efficacy of vaccine and therapeutics (11, 12). Dissecting the evolutionary trajectory of the virus in the

60 population provides important clues to understand the viral replication and virus-host interactions
61 and helps designing effective therapeutics.

62 In this study, we used a NCBI dataset consisting of 6698 high-quality SARS-CoV-2 whole
63 genome sequences with sample collection dates ranging from Dec. 20, 2019 to Jun. 8, 2020. By
64 extensive sequence analysis, we identified the significant and convergent features of the accumulated
65 viral mutations and CpG variations over time. Specifically, in the 29903nt viral genome, four
66 significant mutations, i.e., 5'UTR_c-241-t, NSP3_c-3037-t, NSP12_c-14408-t, and S_a-23403-g,
67 were found to become the dominant variants since early March, 2020, and each of them reached
68 almost 100% of all virus sequences. By global statistical analyzing, we identified 14 mutation sites
69 with significant high rates. In addition, we evaluated the mutation trajectories by each day and every
70 10 days, and notably identified three co-mutation patterns consisting of these 13 sites (among these
71 14 sites) with surprisingly high co-occurrence rates. Moreover, we found the significant decrease of
72 CpG dinucleotides in the viral genome over time, suggesting an evolutionary escape of host innate
73 immunity of CpG (13-15). The dissected evolution trajectory that the four dominant mutations, three
74 significant multi-site co-mutations, and CpG (decrease) mutation contribute to the rapid evolution of
75 SARS-CoV-2 virus in the population, which shapes the viral infectivity and pathogenicity. This
76 study provides valuable clues and frameworks to dissect the viral replication and virus-host
77 interactions for designing effective therapeutics.

78 RESULTS

79 Dominant mutations appeared in SARS-CoV-2 in COVID-19 population over time

80 To explore the mutational landscape of SARS-CoV-2 during virus circulating in the COVID-19
81 population since the outbreak of COVID-19, we aligned 6698 high quality full-length genome
82 sequences across all major regions with viral sample collection dates ranging from Dec. 20, 2019 to
83 Jun. 8, 2020. As the mutation landscape was massive up to date, we identified 82 mutation sites with
84 mutation rate >1% to draw a heatmap (Fig. 1A). As shown in Fig. 1, the Y-axis represents the
85 collection dates of COVID-19 samples, each of which contains 1~225 sequences. According to the
86 first posted viral sequence (NC_045512), there were accumulated mutations during the virus
87 circulating and apparently new mutation sites gradually emerged since the end of Feb. 2020. Noted
88 that several highly mutated sites appeared before Feb. 22, 2020, which was likely due to the limited
89 collected sequences available at that time and accidental random mutations, i.e., a high mutation rate
90 resulted from even one mutation site. Up to now, there were at least four dominant mutations

91 (5'UTR_c-241-t, NSP3_c-3037-t, NSP12_c-14408-t, and S_a-23403-g) (Fig. 1A), where S_
92 a-23403-g mutation resulted in the amino acid change (D614G) that enhances viral infectivity (6, 16),
93 albeit debate exists (10). In particular, each of them covered almost 100% of all virus sequences
94 since the middle May 2020. Focusing on eight mutation sites (5'UTR_c-241-t, NSP3_c-3037-t,
95 NSP12_c-14408-t, S_a-23403-g, ORF3a_g-25563-t, N_g-28881-a, N_g-28882-a, and
96 N_g-28883-c), all the sites sites began to have very high mutation rates since May 2020 (Fig. 1B to
97 1F). Notably, mutations in three adjacent sites in N (N_g-28881-a, N_g-28882-a, and N_g-28883-c)
98 co-occurred (Fig. 1G-I), suggesting a strong selection pressure.

99 **Strong co-occurrent mutations appeared on multiple sites over time**

100 We assessed the mutations of all residues of the SARS-CoV-2 based on the collected genome
101 sequences. The top 34 mutation sites (with mutation rate>2%) were listed in Fig. 2A. Clearly, there
102 were four dominant mutants (S_a-23403-g, NSP12_c-14408-t, 5'UTR_c-241-t, and NSP3_c-3037-t).
103 The three adjacent sites in Nucleocapsid (N) also had considerably high mutation rates (>0.08). By
104 analyzing the top 34 mutations, there were biased mutation patterns, e.g. the ratio of c-to-t (c-t) was
105 more than 44% (Fig. 2B). We then studied the global co-occurrence relationships of the 34 mutations.
106 We found that there were strong co-occurrence site pairs/associations (Fig. 2C) such as the following
107 three multi-site patterns (i) S_a-23403-g, NSP12_c-14408-t, 5'UTR_c-241-t, NSP3_c-3037-t,
108 ORF3a_c-25563-t, and ORF3a_g-25563-t; (ii) ORF8_t-28144-c, NSP4_c-8782-t, NSP14_c-18060-t,
109 NSP13_a-17858-g, and NSP13_c-17747-t; and (iii) N_g-28881-a, N_g-28882-a, and N_g-28883-c.
110 We further quantified the co-occurrence significance (the ratio of co-occurrence mutants to all
111 sequence examples, also called Support (17-19) in data mining) among the top 11 mutation sites
112 (mutation rate>10%) (Fig. 2E to 2I). The Y-axes represented the co-occurrence site pairs. For
113 simplicity, we used gene names instead of their mutation sites in Fig. 2E to 2I. The corresponding
114 mutational positions of genes were given in Fig. 2E in details. From Fig. 2F, every mutational
115 association was significant (>0.10) in all collected genome sequences and almost all associations
116 (except NSP4-NSP13a pair) met strong co-occurrence relationships (>0.60). Figs. 2H-2J show 3-to-6
117 mutation-site (multi-site) co-occurrences. Interestingly, all mutational associations in Fig. 2H to 2J
118 follow significant and strong co-occurrence relationships. The significant co-occurrences of
119 multi-site mutations may suggest that these sites closely interact with each other during the
120 evolution.

121

122 Co-evolution of multi-sites in COVID-19 population

123 The strong co-occurrence relationships of multi-sites in COVID-19 virus suggested
124 co-mutational evolution. We then investigated the co-occurrences of three groups involving 14 sites
125 on a time scale of about per ten days. As shown in Fig. 3A, the mutation rates of the first group
126 containing 6 sites exhibited very similar trends. Strikingly, 4 mutant sites (S_23403, NSP12_14408,
127 5'UTR_241, and NSP3_3037) almost shared a same mutation rate curve. A second group involved 5
128 sites (ORF8_28144, NSP4_8782, NSP14_18060, NSP13_17858, and NSP13_17747) and two
129 mutant sites (ORF8_28144 and NSP4_8782) shared a same mutation rate curve whereas the other
130 three mutant sites (NSP14_18060, NSP13_17858, and NSP13_17747) almost had a same mutation
131 rate curve (Fig. 3B). We analyzed the correlations of the 11 mutant sites of the first and the second
132 groups on a 15 intervals by Pearson Correlation Coefficient (PCC) (20-22), and expressed them by
133 heatmap (Fig. 3C). There were two red (6*6 and 5*5) regions that corresponded to the first and the
134 second mutant groups, respectively. As expected, the top 5 mutation sites, especially the top 4 sites
135 in the first group exhibited a very strong correlation. In the second red (5*5) region corresponded to
136 the second group, there were two subgroups containing a two-site (ORF8_28144 and NSP4_8782)
137 and a triple-site (NSP14_18060, NSP13_17858 and NSP13_17747) exhibited very strong
138 correlations, respectively (Fig. 3C).

139 Unlike the first and the second groups, the third group consisted of three *adjacent* mutation sites
140 in nucleocapsid region (N_28881, N_28882, and N_28883). The mutation rate curve of these sites
141 almost overlapped with each other and the mutation rates of these sites increased over time (Fig. 3D).
142 The correlations of these sites were nearly 1.0 (Fig. 3D), suggesting a strong co-evolution.

143 Based on the relationships of the top two multi-site mutation groups, we illustrated the
144 correlation networks of these mutations. As shown in Fig. 3E, the correlation networks of the first
145 group mutation sites showed a six-pointed star and that of the second group showed a five-pointed
146 star network. The mutation sites exhibited strong correlations with each other (Fig. 3E) and the
147 correlations were further evidenced by the heatmap as show in Fig. 3F.

148 The mutations either changed the amino acid sequence or not. The mutation NSP2_c-1059-t
149 resulted in an amino acid change (T85I) in the NSP2; the mutation NSP12_c-14408-t resulted in an
150 amino acid change in the viral RNA-dependent RNA polymerase NSP12 (P323L) (fig. S1); the
151 mutation S_a-23403-g resulted in an amino acid change in Spike protein (S)(D614G), which
152 enhances viral infectivity (6, 16); the mutation ORF3a_g-25563-t resulted in an amino acid change
153 (Q57H) in Spike protein in ORF3a; the mutation NSP13_c-17747-t and NSP13_a-17858-g resulted

154 in an amino acid changes P504L and Y541C in the NTPase/helicase domain NSP13, respectively;
155 the mutation ORF8_t-28144-c resulted in an amino acid change (L84S) in the ORF8; the mutations
156 N_g-28881-a, N_g-28882-a, N_g-28883-c resulted amino acid changes R203K, R203S and G203R
157 in the nucleocapsid N, respectively (Fig. 3G). In contrast, the mutations NSP14_c-18060-t and
158 NSP3_c-3037-t were synonymous mutations without amino acid changes (Fig. 3G). The mutation
159 5'UTR_c-241-t located in the 5'-non-translated region and did not change the predicted RNA
160 structure of 5'UTR (fig. S2).

161 **Significant decrease of CpG dinucleotide content in 5'UTR in COVID-19 population over time**

162 The CpG content of viral genome is restricted by host intrinsic zinc finger antiviral protein that
163 interacts with CpG rich-region and mediates depletion of foreign viral RNAs (14, 23). Comparing
164 with other coronaviruses, SARS-CoV-2 genome exhibits extreme CpG deficiency (13). However, the
165 evolutionary trajectory of SARS-CoV-2 CpG-content within the same species is still unclear. We
166 investigated the CpG-content changes in SARS-CoV-2 since the outbreak. As shown in Fig. 4A and
167 4B, the CpG dinucleotide content exhibited a decreased trend over time. The CpG-content in each
168 SARS-CoV-2 genome regions varied, with high CpG-contents the 5'UTR, NSP1, E and ORF10
169 regions and low CpG-contents in NSP8, ORF6 regions. Notably, the NSP7 region was free of CpG
170 dinucleotide (Fig. 4C). Comparing with the first posted SARS-CoV-2 genome (NC_045512), in the
171 very recent SARS-CoV-2 genomes, only the CpG-contents of 5'UTR decreased significantly but not
172 the other CpG high content regions NSP1, E and ORF10 (Fig. 4D, E), suggesting a biased evolution
173 pressure on this region.

174 **CONCLUSION**

175 Our comprehensive and massive mutation and correlation analyses identified four dominant
176 mutation sites (5'UTR_c-241-t, NSP3_c-3037-t, NSP12_c-14408-t, and S_a-23403-g) and revealed
177 three significant multi-site co-mutational patterns (S_a-23403-g, NSP12_c-14408-t, 5'UTR_c-241-t,
178 NSP3_c-3037-t, ORF3a_c-25563-t; ORF8_t-28144-c, NSP4_c-8782-t, NSP14_c-18060-t,
179 NSP13_a-17858-g, NSP13_c-17747-t; and N_g-28881-a, N_g-28882-a, N_g-28883-c). Some of the
180 mutations changed the amino acid sequence in the viral RNA-dependent RNA polymerase nsp12
181 (P323L), Spike protein (S) (D614G), ORF3a (Q57H), NTPase/helicase domain nsp13 (P504L,
182 Y541C), ORF8 (L84S) and nucleocapsid N (R203K, R203S, G203R), or did not change the amino
183 acid sequence (NSP14_c-18060-t and NSP3_c-3037-t), which may affect viral replication or
184 virus-host interaction. Other mutations were synonymous mutations without amino acid changes (Fig.

185 3G). And mutations located in the 5'-non-translated region (5'UTR_c-241-t) that did not change the
186 predicted RNA structure. Moreover, we found gradual but significant decrease of CpG-content in
187 5'UTR region over time, which suggested the 5'UTR region as a potential ZAP target. Taken together,
188 our study provides valuable clues and frameworks to dissect the viral replication and virus-host
189 interactions for designing effective therapeutics.

190
191 **Acknowledgments:** We thank associate prof. T.Z. and Dr. S.T.H. for useful comments on the
192 manuscript.

193 **Funding:** This work is supported by the National Key Research and Development Program of China
194 (2017YFA0505500 to L.N.C., 2017YFC0909502 to J.S.Z.); the Strategic Priority Research Program
195 of the Chinese Academy of Sciences (XDB38040400 to L.N.C.); National Science Foundation of
196 China (31771476 and 31930022 to L.N.C, 61602460 and 11701379 to J.S.Z.); Shanghai Municipal
197 Science and Technology Major Project (2017SHZDZX01 to L.N.C.); National Science and
198 Technology Major Project of China (2017ZX10103009 to Z.G.Y.); Emergency Project of Shanghai
199 Science and Technology Committee (20411950103 to Z.G.Y.); Development programs for
200 COVID-19 of Shanghai Science and Technology Commission (20431900401 to Z.G.Y.); National
201 Postdoctoral Program for Innovative Talent (BX20180331 to J.Y.K.); and China Postdoctoral
202 Science Foundation (2018M642018 to J.Y.K.).

203 **Author contributions:** L.N.C. and J.S.Z. designed the study. J.S.Z. and Z.G.Y. designed the
204 experiments. J.S.Z. analyzed data. J.S.Z., Z.G.Y., and J.Y.K. designed the figures. H.B.H. checked
205 the experiments. J.S.Z. wrote the manuscript. Z.G.Y., J.Y.K., M.F.L., L.L., and Y.Q.H polished the
206 manuscript.

207 **Data and materials availability:** All data are available in the main text or supplementary materials
208 or from the corresponding author upon request.

209 **Conflict of interest**

210 The authors declare no conflict of interest.

211

212 **Supplementary Materials:**

213 Figs. S1 and S2

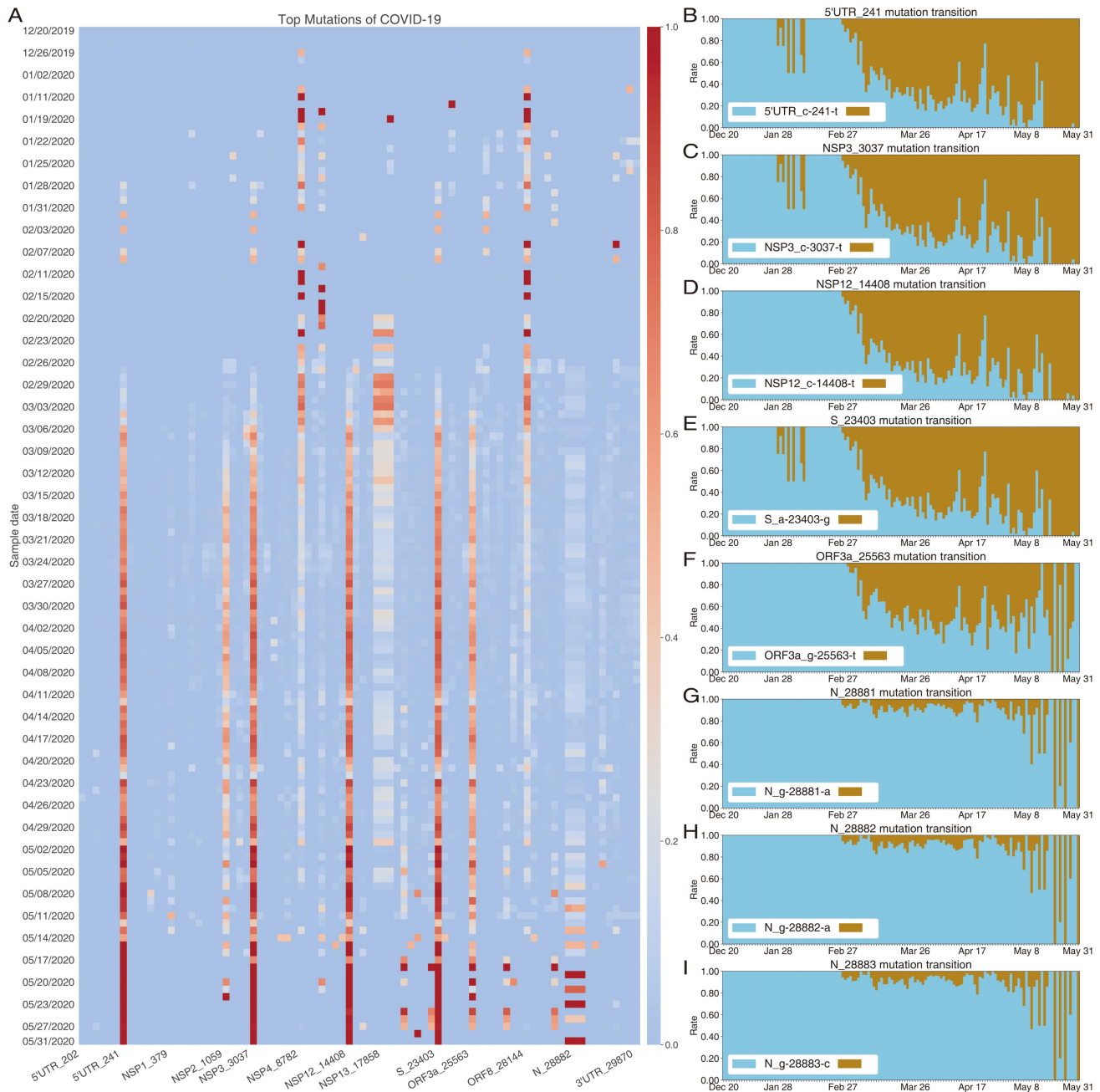
214

215 **References and Notes:**

- 216 1. D. Wrapp *et al.*, Cryo-EM structure of the 2019-nCoV spike in the prefusion conformation. *Science* **367**,
217 1260-1263 (2020).
218 2. M. Chinazzi *et al.*, The effect of travel restrictions on the spread of the 2019 novel coronavirus (COVID-19)

- 219 outbreak. *Science*, 395-400 (2020).
- 220 3. N. Zhu *et al.*, A Novel Coronavirus from Patients with Pneumonia in China, 2019. *The New England journal of*
221 *medicine* **382**, 727-733 (2020).
- 222 4. P. Zhou *et al.*, A pneumonia outbreak associated with a new coronavirus of probable bat origin. *Nature* **579**,
223 270-273 (2020).
- 224 5. L. Chen *et al.*, RNA based mNGS approach identifies a novel human coronavirus from two individual
225 pneumonia cases in 2019 Wuhan outbreak. *Emerging microbes & infections* **9**, 313-319 (2020).
- 226 6. B. Korber *et al.*, Tracking changes in SARS-CoV-2 Spike: evidence that D614G increases infectivity of the
227 COVID-19 virus. *Cell*, (2020).
- 228 7. E. C. Smith, H. Blanc, M. Vignuzzi, M. R. Denison, Coronaviruses Lacking Exoribonuclease Activity Are
229 Susceptible to Lethal Mutagenesis: Evidence for Proofreading and Potential Therapeutics. *Plos Pathog* **9**,
230 (2013).
- 231 8. M. Sevajol, L. Subissi, E. Decroly, B. Canard, I. Imbert, Insights into RNA synthesis, capping, and proofreading
232 mechanisms of SARS-coronavirus. *Virus Res* **194**, 90-99 (2014).
- 233 9. M. Vignuzzi, J. K. Stone, J. J. Arnold, C. E. Cameron, R. J. N. Andino, Quasispecies diversity determines
234 pathogenesis through cooperative interactions in a viral population. **439**, 344-348 (2006).
- 235 10. N. D. Grubaugh, W. P. Hanage, A. L. Rasmussen, Making sense of mutation: what D614G means for the
236 COVID-19 pandemic remains unclear. *Cell*, (2020).
- 237 11. B. A. *et al.*, Antibody cocktail to SARS-CoV-2 spike protein prevents rapid mutational escape seen with
238 individual antibodies. *Science*, (2020).
- 239 12. J. Hansen *et al.*, Studies in humanized mice and convalescent humans yield a SARS-CoV-2 antibody cocktail.
240 *Science*, (2020).
- 241 13. X. Xia, Extreme Genomic CpG Deficiency in SARS-CoV-2 and Evasion of Host Antiviral Defense. *Molecular*
242 *Biology and Evolution*, (2020).
- 243 14. J. L. Meagher *et al.*, Structure of the zinc-finger antiviral protein in complex with RNA reveals a mechanism for
244 selective targeting of CG-rich viral sequences. *P Natl Acad Sci USA* **116**, 24303-24309 (2019).
- 245 15. M. Ficarelli *et al.*, KHNYN is essential for the zinc finger antiviral protein (ZAP) to restrict HIV-1 containing
246 clustered CpG dinucleotides. *Elife* **8**, (2019).
- 247 16. Q. Li *et al.*, The impact of mutations in SARS-CoV-2 spike on viral infectivity and antigenicity. *Cell*, (2020).
- 248 17. J. S. Zhang, Y. L. Wang, C. Zhang, Y. Y. Shi, Mining Contiguous Sequential Generators in Biological Sequences.
249 *Ieee Acm T Comput Bi* **13**, 855-867 (2016).
- 250 18. J. S. Zhang, Y. L. Wang, D. Y. Yang, CCSpan: Mining closed contiguous sequential patterns. *Knowl-Based Syst*
251 **89**, 1-13 (2015).
- 252 19. J. S. Zhang *et al.*, Efficient Mining Multi-Mers in a Variety of Biological Sequences. *Ieee Acm T Comput Bi* **17**,
253 949-958 (2020).
- 254 20. W. Wiedermann, M. Hagmann, Asymmetric properties of the Pearson correlation coefficient: Correlation as the
255 negative association between linear regression residuals. *Commun Stat-Theor M* **45**, 6263-6283 (2016).
- 256 21. Y. S. Mu, X. D. Liu, L. D. Wang, A Pearson's correlation coefficient based decision tree and its parallel
257 implementation. *Inform Sciences* **435**, 40-58 (2018).
- 258 22. G. Y. Mao, On high-dimensional tests for mutual independence based on Pearson's correlation coefficient.
259 *Commun Stat-Theor M* **49**, 3572-3584 (2020).
- 260 23. V. Odon *et al.*, The role of ZAP and OAS3/RNaseL pathways in the attenuation of an RNA virus with elevated
261 frequencies of CpG and UpA dinucleotides. *Nucleic Acids Res* **47**, 8061-8083 (2019).
- 262

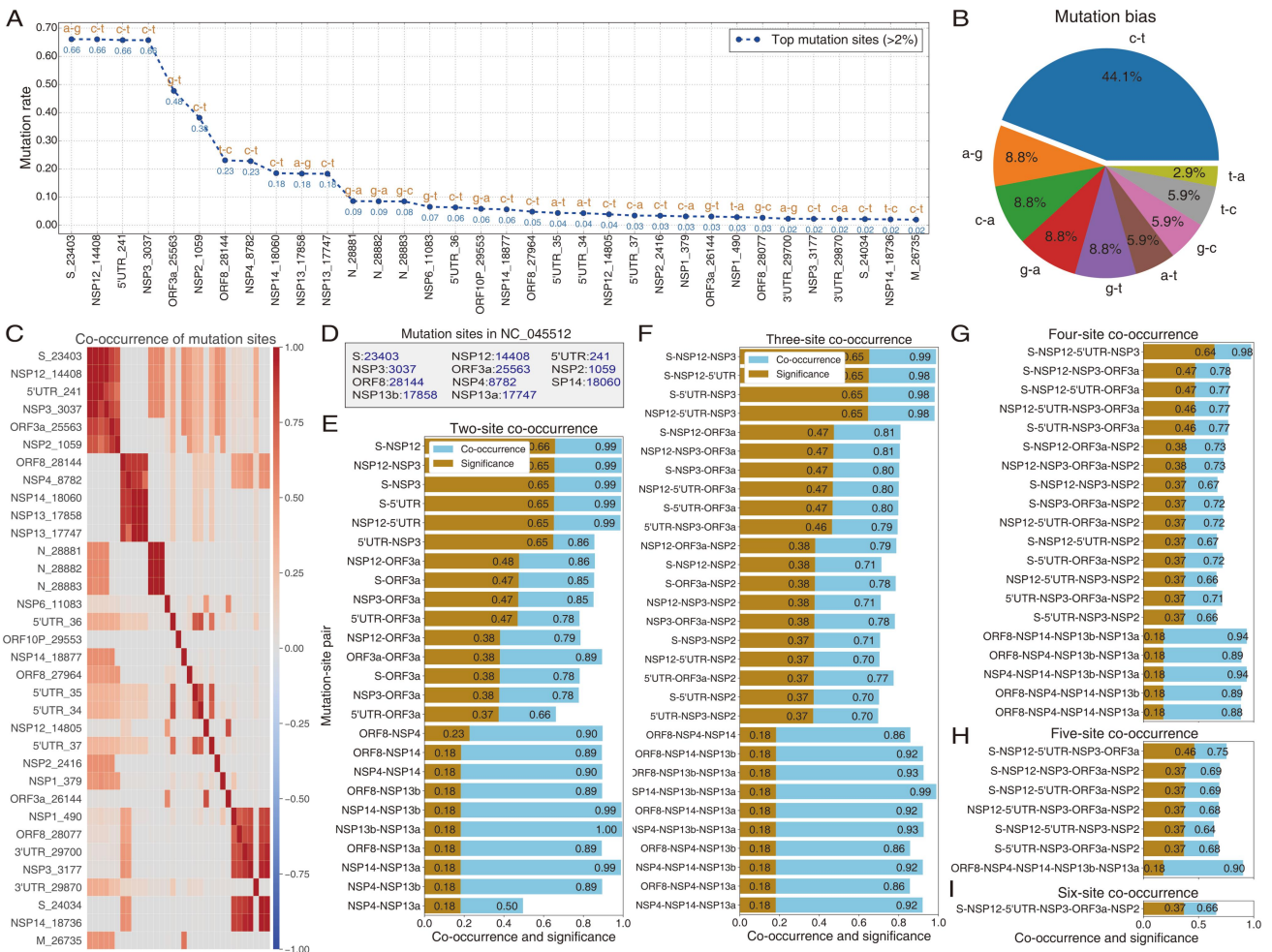
263 **Figure 1**



264

265 **Fig. 1. Ongoing and dominant mutations of SARS-CoV-2 over time.** (A) The global mutational
266 landscape of top 1% mutation rates from Dec. 20, 2019 to May 31, 2020. The mutation of
267 SARS-CoV-2 is clearly ongoing and yet with a rapid rate. (B)- (I) show the mutation transitions of 8
268 sites from Dec. 20, 2019 to May 31, 2020. The sky-blue represents the rates of the original nucleic
269 acids in reference sequence, and the dark-golden the rates of the mutant nucleic acids. Note that the 5
270 mutations (subfigures B to F), i.e., 5'UTR_c-241-t, NSP3_c-3037-t, NSP12_c-14408-t, and
271 S_a-23403-g, and ORF3a_c-25563-t, especially the top 4 ones have clearly become the dominant
272 mutants. The three adjacent mutations (subfigures G to I, N_g-28881-a, N_g-28882-a, and
273 N_g-28883-c) increase daily on the whole.

274 **Figure 2**

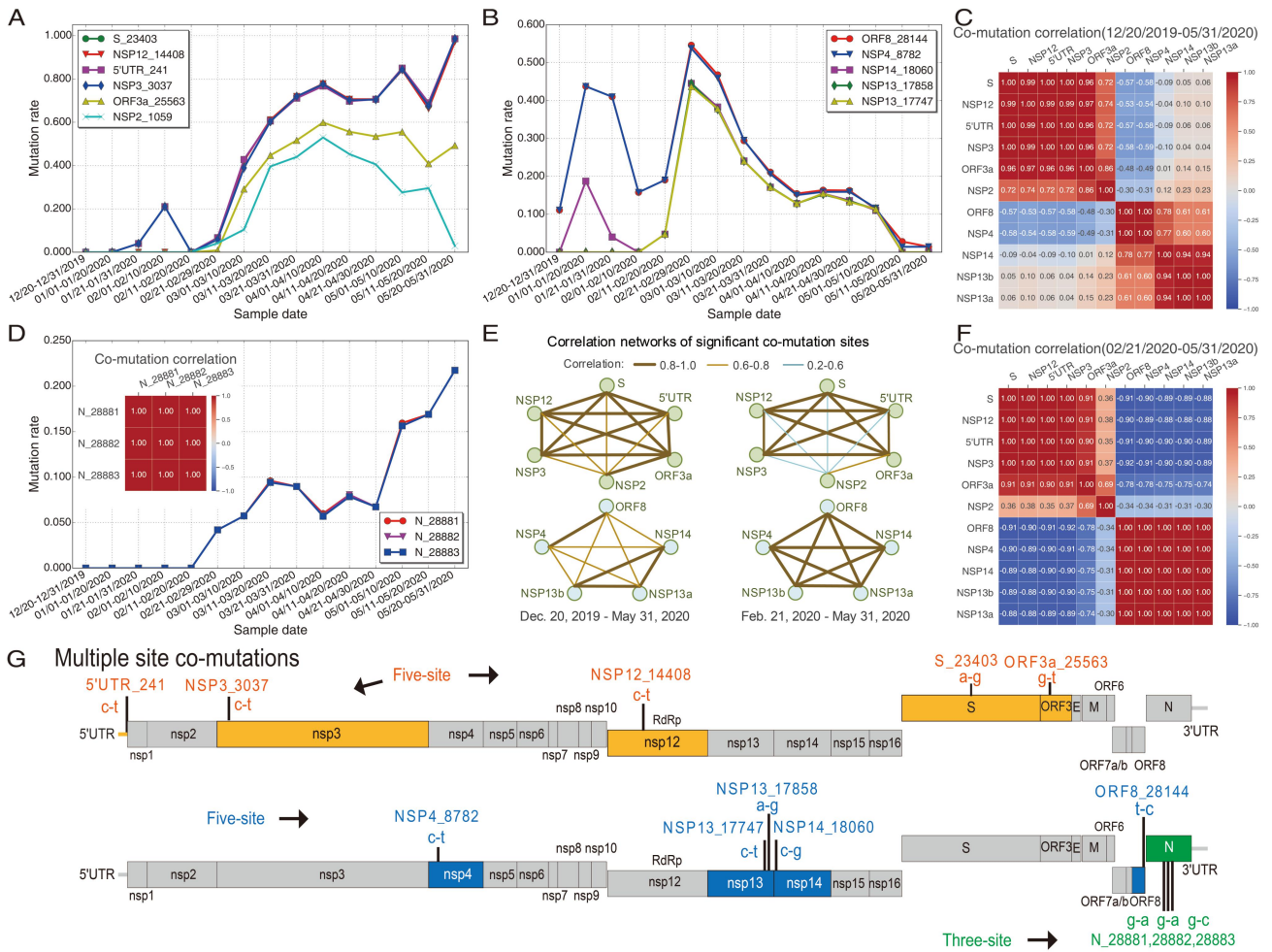


275

276 **Fig. 2. Co-occurrence mutations of SARS-CoV-2 over time. (A)** Top 2% mutation sites involve 34
 277 sties. The top 11 sites hold significant high mutation rates. **(B)** Mutation bias of the top 34 mutants.
 278 **(C)** Co-occurrence heatmap of top 2% mutation sites in order of mutational significance. The X-axis
 279 and Y-axis share the same tick labels (mutation sites along with their positions in reference sequence
 280 NC_045512) as shown in Y-axis. The top 14 sites were clustered into three high co-occurrence
 281 groups. **(D)** Query table of top 11 mutation sites. We used gene names instead of their mutation sites
 282 in subfigures (E) to (I), for simplicity. **(E)** Co-occurrences of 2 mutation sites. Each tick label of
 283 Y-axis, like S-NSP12, represents the mutation-site association/pattern of gene S (position 23403) and
 284 NSP12 (position 14408). The detailed positions of these mutations are shown in subfigures (E). **(F)**
 285 Co-occurrences of 3 mutation sites. **(G)** Co-occurrences of 4 mutation sites. **(H)** Co-occurrences of 5
 286 mutation sites. **(I)** Co-occurrences of 6 mutation sites.

287

288 **Figure 3**



289

290 **Fig. 3. Co-evolution of multi-sites in COVID-19 population over time. (A)** Mutation trends of top
 291 6 sites in Fig. 2A from Dec. 20, 2019 to May 31, 2020. **(B)** Mutation trends from the 7th to 11th sites
 292 in Fig. 2A. **(C)** Global co-mutational heatmap of top 11 in Fig. 2A sites since Dec. 20, 2019. **(D)**
 293 Mutation trends and co-mutational heatmap of three adjacent Nucleocapsid sites (from 12th to 14th
 294 sites in Fig.2A). **(E)** Correlation networks of significant co-mutation sites. The left two sites show
 295 the correlation relationships from Dec. 20, 2019 to May 31, 2020, and the right two from Feb. 21,
 296 2020 to May 31, 2020. **(F)** Co-mutational heatmap of top 11 sites with sample collection dates from
 297 Feb. 21, 2020 to May 31, 2020. **(G)** The mutational positions of multi-site co-mutations consisting of
 298 5-, 5- and 3-site co-mutations.

299

300

301

302

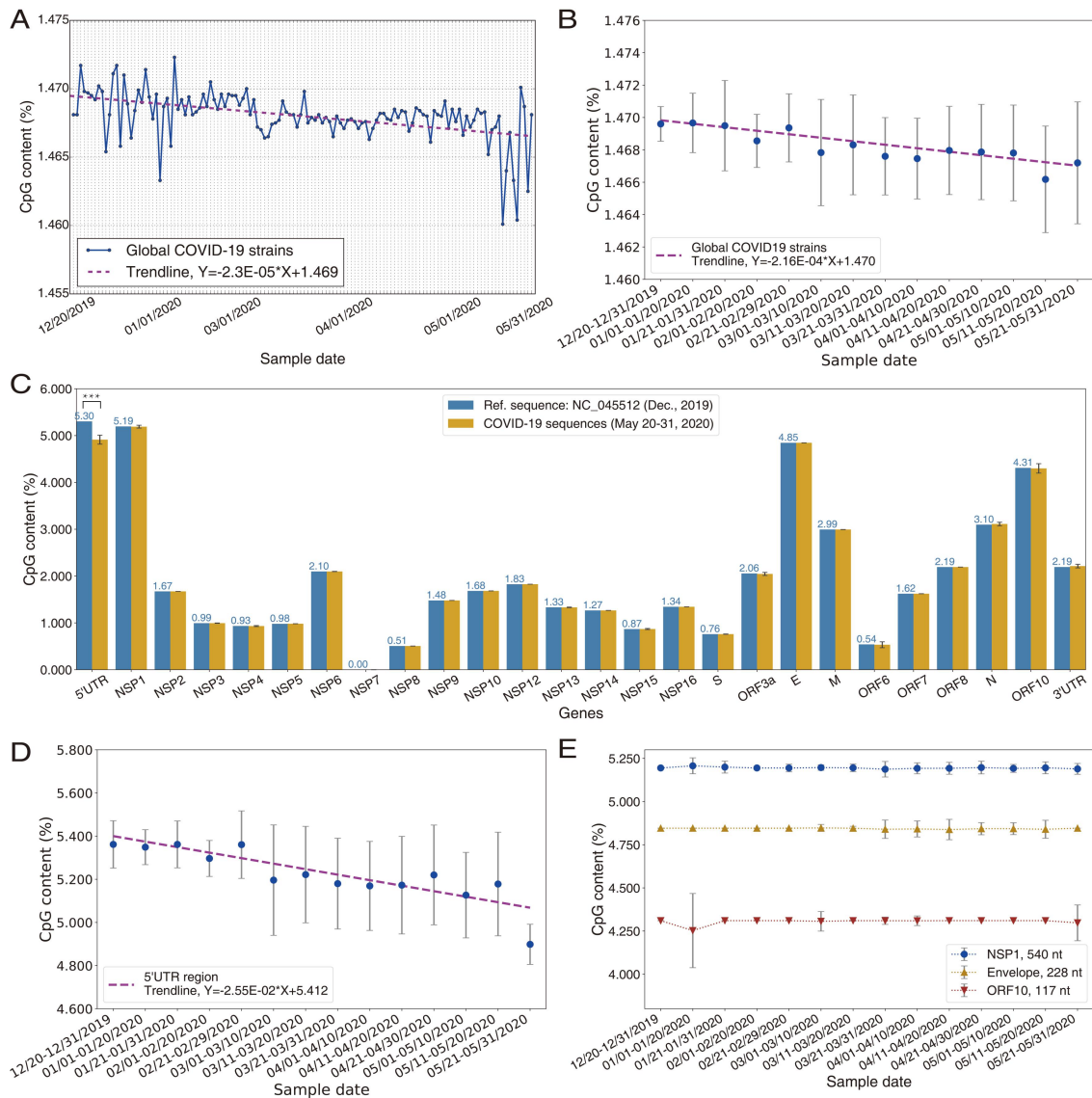
303

304

305

306

307 **Figure 4**



308

309 **Fig. 4. CpG-content decrease in COVID-19 population over time. (A)** Trend of CpG decrease of
310 COVID-19 genome per day. **(B)** Trend of CpG decrease of COVID-19 genome by intervals of about
311 10 days. **(C)** CpG changes of all genes. The significant decrease of only 5'UTR region indicates that
312 5'UTR is the potential target gene of ZAP. **(D)** Trend of CpG decrease of 5'UTR region by about 10
313 days. **(E)** CpG change trends of NSP1, Envelope, and ORF10.

314