

1 **Human cell-dependent, directional, time-dependent changes in the mono- and**  
2 **oligonucleotide compositions of SARS-CoV-2 genomes**

3

4 Yuki Iwasaki<sup>1</sup>, Takashi Abe<sup>2</sup>, Toshimichi Ikemura<sup>1</sup>

5 1. Department of Bioscience, Nagahama Institute of Bio-Science and Technology.

6 Shiga, Japan

7 2. Graduate School of Science and Technology, Niigata University, Niigata, Japan

8

9 **Abstract**

10 **Background**

11 When a virus that has grown in a nonhuman host starts an epidemic in the human  
12 population, human cells may not provide growth conditions ideal for the virus.

13 Therefore, the invasion of severe acute respiratory syndrome coronavirus-2 (SARS-  
14 CoV-2), which is usually prevalent in the bat population, into the human population is  
15 thought to have necessitated changes in the viral genome for efficient growth in the new  
16 environment. In the present study, to understand host-dependent changes in coronavirus  
17 genomes, we focused on the mono- and oligonucleotide compositions of SARS-CoV-2  
18 genomes and investigated how these compositions changed time-dependently in the  
19 human cellular environment. We also compared the oligonucleotide compositions of  
20 SARS-CoV-2 and other coronaviruses prevalent in humans or bats to investigate the  
21 causes of changes in the host environment.

22 **Results**

23 Time-series analyses of changes in the nucleotide compositions of SARS-CoV-2  
24 genomes revealed a group of mono- and oligonucleotides whose compositions changed  
25 in a common direction for all clades, even though viruses belonging to different clades  
26 should evolve independently. Interestingly, the compositions of these oligonucleotides

27 changed towards those of coronaviruses that have been prevalent in humans for a long  
28 period and away from those of bat coronaviruses.

## 29 **Conclusions**

30 Clade-independent, time-dependent changes are thought to have biological significance  
31 and should relate to viral adaptation to a new host environment, providing important  
32 clues for understanding viral host adaptation mechanisms.

33

## 34 **Keyword**

35 “COVID-19”, “SARS-CoV-2”, “Oligonucleotide composition”, “Time-series analysis”,  
36 “Big data”, “Zoonotic virus”, “RNA virus”, “Viral adaptation”, “Coronavirus”

37

## 38 **Background**

39 Severe acute respiratory syndrome coronavirus-2 (SARS-CoV-2), an RNA virus  
40 belonging to the betacoronavirus genus, began to spread in the human population in  
41 2019. This viral strain is believed to have been originally prevalent in bats and  
42 transferred to the human population through intermediate hosts [1]. Viral growth  
43 requires a wide variety of host factors (nucleotide pools, proteins, RNA, etc.) and  
44 should evade the diverse antiviral mechanisms of host cells (antibodies, killer T cells,  
45 interferon, RNA interference, etc.) [2-4]. Since ancestral SARS-CoV-2 strains are  
46 thought to be endemic in bats, they should be well adapted to their host environment;  
47 when the virus invades the human population, human cells may not provide growth  
48 conditions ideal for the virus. For efficient growth and rapid spread of the infection,  
49 changes in the viral genome should be required. Analyses of time-dependent changes in  
50 SARS-CoV-2 in the human population can be used to characterize how and why viral  
51 genomes change to adapt to a new host environment.

52         Due to the great threat of COVID-19 and remarkable development of  
53 sequencing technology, a massive number of SARS-CoV-2 genome sequences are

54 available in databases, even though the epidemic has lasted for approximately 10  
55 months. These sequence data have provided a wide range of insights into SARS-CoV-2  
56 [5,6]. Phylogenetic methods based on sequence alignment have been widely used in  
57 molecular evolution studies [7,8], and these methods are well refined and essential for  
58 studying phylogenetic relationships between different viral species and variations in the  
59 same viral species at the single-nucleotide level. However, when dealing with a massive  
60 number of genome sequences, methods based on sequence alignment become  
61 problematic because they require a large amount of computational resources.

62         We have continued to develop sequence alignment-free methods focused on  
63 the oligonucleotide compositions of genome sequences [9-12]. Notably, oligonucleotide  
64 composition varies widely among species, including viruses, and is designated as  
65 genome signatures [13]. These compositions can be treated as numerical data, and a  
66 massive amount of sequence data can easily be subjected to various statistical analyses.  
67 Furthermore, even genomic fragments without orthologous and/or paralogous pairs can  
68 be compared [11,12,14-17]. Specifically, our previous work on influenza A-type virus  
69 genomes found that the oligonucleotide compositions of the viral genomes differed  
70 between hosts (e.g., humans and birds), even for viruses within the same subtype (e.g.,  
71 H1N1 and H3N2 of type A) [11,12,14]; we also examined changes in the  
72 oligonucleotide compositions of influenza H1N1/09, which have been epidemic in  
73 humans beginning in 2009, and found that their compositions changed to approach  
74 those of the seasonal flu strains H1N1 and H3N2 [11]. Furthermore, although epidemics  
75 of the H1N1 and H3N2 strains began several decades apart, these strains showed highly  
76 similar chronological changes from the start of these epidemics. These evolutionary yet  
77 reproducible changes suggest that mutations to adapt to a new host environment  
78 inevitably accumulate when the host species of a virus changes, and these changes can  
79 be efficiently detected by analyzing oligonucleotide compositions.

80           Several groups, including ours, have examined changes in SARS-CoV-2  
81 genomes during the early stages of the SARS-CoV-2 epidemic and found clear  
82 directional changes in a group of mono- and oligonucleotides detectable on even a  
83 monthly basis [15,18,19]. These directional changes will allow us to predict changes in  
84 the near future. Notably, near-future prediction and verification should be the most  
85 direct ways to test the reliability of the obtained results, models and ideas (e.g., those  
86 discovered for influenza viruses), providing a new paradigm for molecular evolutionary  
87 studies. In this context, the present study analyzed the genome sequences of over  
88 seventy thousand SARS-CoV-2 strains isolated from December 2019 to September  
89 2020.

90

91

## 92 **Results**

### 93 **Directional changes in the mononucleotide compositions (%) of SARS-CoV-2**

94 For fast-evolving RNA viruses, diversity within the viral population arises rapidly as  
95 the epidemic progresses and subpopulation structure forms; the GISAID consortium has  
96 defined at least seven main clades (G, GH, GR, L, V, S and Others). Notably, the  
97 elementary processes of molecular evolution are based on random mutations, and  
98 strains belonging to different clades are thought to have evolved independently.  
99 Therefore, the observation of highly similar time-dependent changes independent of  
100 clade has certain biological meanings and may be inevitable for efficient growth in  
101 human cells. From this perspective, we first examined time-dependent changes in the  
102 mononucleotide compositions (%) of SARS-CoV-2 strains isolated from December  
103 2019 to September 2020.

104           Among the seven clades (G, GH, GR, L, V, S and Others) reported by the  
105 GISAID consortium, we used six clades (G, GH, GR, L, V and S), excluding Others, in  
106 the analysis. For the time-series analysis, we calculated the average mononucleotide

107 compositions (%) of the genomes in each clade collected monthly; in Fig. 1A, the  
108 mononucleotide composition of each clade is shown as a colored line, while that for the  
109 monthly collected genomes belonging to all clades is shown as a dashed line.

110         Regardless of clade, the composition of C decreased, while that of U increased  
111 in a time-dependent manner, but the changes in A and G composition were less clear  
112 (Fig. 1A). Correlation coefficients between the mononucleotide composition and month  
113 from the start of the epidemic showed a high negative correlation for C and a high  
114 positive correlation for U for all clades, but there was no clear directionality for A and  
115 G (Fig. 1A and Tables 1, 2). These results indicate that the mononucleotide composition  
116 of this virus may be prone to biased mutations that reduce C and increase U or the  
117 mutated strains tend to be more favorable for growth in human cells.

118

### 119 **Directional changes in short oligonucleotide compositions**

120 Oligonucleotides are known to act as functional motifs, such as binding sites for a wide  
121 variety of proteins and target sites for RNA modifications. Therefore, directional  
122 changes in some oligonucleotides independent of clade may relate to certain processes  
123 for adaptation to the new host environment. Our previous work on influenza A viruses  
124 found that their oligonucleotide compositions varied among prevalent hosts [11,12];  
125 notably, although influenza virus isolated from humans tended to prefer A and U (but  
126 not G and C) more than viruses isolated from birds, the human viruses showed a  
127 preference for GGCG and GGGG, which are G- or C-rich. Importantly, there are  
128 various examples of oligonucleotides whose changes in composition cannot be  
129 explained by changes in mononucleotide composition alone, and these changes may  
130 relate to the molecular mechanisms of viral adaptation to a new host.

131         From this perspective, we next analyzed time-dependent changes in di- and  
132 trinucleotide compositions and found that a group of di- and trinucleotides showed a  
133 highly positive or negative correlation (Figs. 1B, S1 and Tables 1, 2). Interestingly, a

134 group of A- or G-rich oligonucleotides, such as GAG and GGA, showed a high positive  
135 correlation independent of clade, which was not expected from the changes in  
136 mononucleotide compositions alone. To confirm the extent of these changes, we also  
137 calculated the fold change in composition for the first isolated month and the last  
138 examined month (Fig. 2) and found clear increases and decreases in mono- and  
139 oligonucleotide compositions common among the six clades, which supports the result  
140 presented in Fig. 1 and Tables 1 and 2.

141

#### 142 **Changes towards the sequences of other coronaviruses prevalent in humans**

143 In a previous study of SARS-CoV-2 [16], we analyzed mono- and dinucleotide  
144 compositions for the first four epidemic months without separating the sequences by  
145 clade. Notably, the directional changes shown in Figs. 1 and 2 and Tables 1 and 2 were  
146 absolutely consistent with the previous results, even when the six clades were separately  
147 analyzed. In the previous study, time-series analysis of ebolavirus at the beginning of  
148 the epidemic in West Africa in 2014 also showed directional changes in a group of  
149 mono- and dinucleotide compositions, but these directional increases/decreases tended  
150 to slow approximately 10 months after the start of the epidemic. The increase/decrease  
151 trend for SARS-CoV-2 is far from slowing after 10 months, and the next important  
152 questions are how long these directional changes in this virus will last and whether there  
153 are possible goals to these changes.

154 To conduct this near-future prediction, the following information concerning  
155 influenza viruses should be useful. As mentioned before, mono- and oligonucleotide  
156 compositions in influenza H1N1/09 changed towards those of seasonal influenza strains  
157 such as the H1N1 and H3N2 subtypes [11]. Furthermore, all the human subtypes  
158 showed directional changes away from the compositions of all avian influenza A  
159 subtypes and closer to those of the human influenza B type, which has been prevalent  
160 only in humans [14]. If we assume that changes similar to those in the influenza virus

161 will occur, the mono- and oligonucleotide compositions of interest for SARS-CoV-2 are  
162 expected to change towards those of other coronaviruses that have been prevalent in  
163 humans and away from those of coronaviruses prevalent in bats. To test this hypothesis,  
164 we analyzed the following coronaviruses: 238 human-CoV strains (alphacoronaviruses  
165 229E and NL63; betacoronaviruses HKU1 and OC43) and 166 bat-CoV strains  
166 (alphacoronaviruses and betacoronaviruses, including the SARS virus).

167 As shown in Fig. 3A, we compared the mononucleotide compositions of  
168 SARS-CoV-2 with those of the human- and bat-CoV strains; the data for bat SARS  
169 among bat-CoV strains, which is thought to be the original strain that caused the current  
170 COVID-19 pandemic, are marked in pink. Interestingly, concerning the human- and  
171 bat-CoV strains, differences in mononucleotide composition were more pronounced  
172 between hosts than between the alpha and beta lineages, and the levels for all six clades  
173 of SARS-CoV-2 were between those for the two hosts. Fig. 3B shows the results of di-  
174 and trinucleotides, for which the directional, time-dependent changes were primarily  
175 common among the six clades. The increases and decreases in nucleotide composition  
176 observed for SARS-CoV-2 in Figs. 1 and 2 are indicated by hollow up and down  
177 arrows, respectively. Interestingly, all changes of interest tended to move away from the  
178 compositions of bat SARS and approach those of human-CoV, supporting the view that  
179 the directional changes of interest have biological significance and are possibly  
180 inevitable, as observed for influenza viruses. Assuming that approaching the levels in  
181 human-CoV strains is the hypothetical goal of the directional change of SARS-CoV-2,  
182 the current compositions are far from this hypothetical goal (Fig. 3); therefore, we  
183 predict that directional changes of interest will continue in the near future.

184 Then, assuming that the average value for all human-CoV strains is a  
185 hypothetical goal, we investigated how SARS-CoV-2 has approached this possible goal.  
186 Specifically, we calculated the square of the difference between the composition of each  
187 nucleotide in SARS-CoV-2 and the average value for human-CoV strains and plotted

188 the values of the difference according to the elapsed month for each nucleotide.  
189 Changes in the compositions of both C and U clearly reduced this difference, as the  
190 compositions of these nucleotides approached the hypothetical goal (Fig. 4A); their  
191 linear reduction supports the prediction that directional changes in the composition of C  
192 and U will continue for the foreseeable future. In contrast, A and G did not show  
193 directional changes in composition, which is most likely due to the absence of clear  
194 differences in the A and G compositions of human- and bat-CoV, i.e., there is no  
195 possible target (Fig. 3A). Fig. 4B shows examples of di- and trinucleotides whose  
196 compositions have moved towards the hypothetical goal, but Fig. 4C shows a few  
197 exceptional nucleotides whose compositions have not changed towards the hypothetical  
198 goal but have changed with a common directionality among the six clades. In Fig. 4D,  
199 correlation coefficients between the above difference and the elapsed month are  
200 presented. Most nucleotides of interest showed a negative coefficient (i.e., a directional  
201 change towards human-CoV), but three oligonucleotides, GG, AGC and CAU, showed  
202 positive coefficients indicating an increase in the difference (i.e., moving away from the  
203 human-CoV level). For these opposing directional changes, certain causes specific to  
204 SARS-CoV-2 may be assumed.

205

### 206 **Motifs for RNA-binding proteins**

207 Next, we considered the mechanisms that move oligonucleotide compositions away  
208 from those of bat coronaviruses and closer to those of human coronaviruses. Certain  
209 human cellular factors involved in viral growth may be candidates in such mechanisms.  
210 When considering possible protein factors, oligonucleotides longer than trinucleotides  
211 should be a focus. As an attempt, we here focused on host RNA-binding proteins  
212 because their binding to hepatitis C virus is known to be involved in the growth of this  
213 RNA virus [20]. We thus searched for motifs for human RNA-binding proteins in  
214 coronavirus genomes (see Methods section) and found multiple loci with binding motifs



215 for each protein. Table 3 (and Table S10) lists the motifs for which a directional time-  
216 dependent change was primarily common among six clades. Table 3 and Fig. 5A show  
217 that only ELAVL1 showed a positive correlation, but the other nine proteins in Table 3  
218 showed a negative correlation for almost all clades; the results for other motifs are  
219 presented in Table S12.

220 We next compared the numbers of these motifs in SARS-CoV-2 with the  
221 numbers of human- and bat-CoV motifs (Fig. 5B). Of the ten proteins shown in Table 3,  
222 the only elevated motif, that for ELAVL1 binding, was found in a significantly higher  
223 number of loci in human-CoV than in bat-CoV, but motifs for PCBP2 and SRSF1  
224 binding, which tended to decrease (Table 3), were found in significantly fewer loci in  
225 human-CoV. These observations appear to be consistent with the features found in the  
226 mono-, di- and trinucleotide compositions of interest. However, unlike these changes,  
227 there was significant diversity within even a single clade, which appears to be greater  
228 than the differences between hosts, with the possible exception of ELAVL1. In regard  
229 to long oligonucleotides, they should carry out a variety of functions, and mutations that  
230 accumulate in their functional motifs may have complex effects on the presence of  
231 functional motif sequences, so an analysis from a new perspective appears to become  
232 important.

233

## 234 **Discussion**

235 We first discuss possible molecular mechanisms related to time-dependent directional  
236 changes in mononucleotide composition. Fig. 1A shows that the frequency of C tended  
237 to decrease in SARS-CoV-2, while that of U tended to increase. Since a similar change  
238 was previously found for MERS and all A-type influenza subtypes [12,14], these  
239 changes may have biological significance for a wide range of RNA viruses that invade  
240 from nonhuman hosts. One possible mechanism is the host RNA-editing function;  
241 Simmonds (2020) proposed that the C→U hypermutation in SARS-CoV-2 may be due

242 to the influence of APOBEC family proteins in humans [19]. APOBEC is an antiviral  
243 protein in various animal species, including humans, that can convert C to U by the  
244 deacetylation of C [21-23]. Such RNA editing is also known to act as a defense  
245 mechanism against various viruses, including retroviruses [24]. The APOBEC gene  
246 family has generated various paralogs during mammalian evolution, with seven known  
247 APOBEC genes in humans and ten in bat families [25-27]. The prevalence C→U  
248 change in SARS-CoV-2 upon transfer of its host environment from bats to humans  
249 suggests that these changes may be due to human-specific APOBEC genes.

250         We next discuss changes in short oligonucleotides. Directional changes in  
251 some oligonucleotides, such as GAG and GGA, cannot be explained by APOBEC-  
252 induced C→U mutations alone. Although the evidence is weak, these oligonucleotides  
253 are part of the binding motifs of several RNA-binding proteins, such as SRSF1 and  
254 PCBP2 (Table S9); the number of loci for these motifs has decreased independently of  
255 clade. In contrast, the number of motif loci for only ELAVL1 among the ten proteins  
256 listed in Table 3 has increased independently of clade. As an RNA-binding protein that  
257 binds A- or U-rich elements, ELAVL1 binding to mRNA is known to contribute to  
258 RNA stability [28, 29]; SARS-CoV-2 and human-CoV, which are prevalent in humans,  
259 may contain increased binding motifs for ELAVL1 for efficient growth in the human  
260 cellular environment. However, for further analysis, information on RNA-binding  
261 proteins in bat cells is needed.

262

## 263 **Conclusions**

264 In the present study, we found that the compositions of a group of mono- and  
265 oligonucleotide in SARS-CoV-2 genomes have changed in a host cell-dependent  
266 manner. This is totally consistent to our previous finding for influenza A and B viruses  
267 [11,12,14], supporting the previous prediction that the host-dependent directional  
268 changes of various mono- and oligonucleotides should inevitably occur in zoonotic

269 RNA viruses that have invaded from nonhuman hosts. Phylogenetic methods based on  
270 sequence alignment [7,8] are well refined and undoubtedly essential for studying the  
271 phylogenetic relationships between viruses. The present alignment-free method to  
272 analyze mono- and oligonucleotide compositions can also serve as a powerful tool for  
273 molecular evolutionary studies of viruses, revealing directional changes in viruses and  
274 predicting the possible goals of these changes.

275

276

## 277 **Methods**

### 278 **SARS-CoV-2 genome sequences**

279 Human SARS-CoV-2 genome sequences were downloaded from the GISAID database  
280 (<https://www.gisaid.org/>); sequences that were complete, showed high coverage and had  
281 been isolated from humans were downloaded on Sep 17, 2020. Among the acquired  
282 sequences, strains with an unknown isolation month were excluded from the analysis,  
283 and the polyA tail was removed. A list of all 72,314 strains used is provided in Table  
284 S1.

285

### 286 **Genome sequences of coronaviruses prevalent in humans or bats**

287 The complete sequences of two types of human coronavirus (human-CoV) strains,  
288 alphacoronaviruses (27 229E and 55 NL63 strains) and betacoronaviruses (18 HKU1  
289 and 138 OC43 strains), were obtained from the NCBI virus database  
290 (<https://www.ncbi.nlm.nih.gov/labs/virus/>). The complete genome sequences of two  
291 types of bat coronavirus (bat-CoV) strains, alphacoronaviruses (87 strains) and  
292 betacoronaviruses (79 strains, including 34 SARS-CoV), isolated from three types of  
293 bats (Chiroptera, Vespertilionidae and Rhinolophidae) were obtained from the NCBI  
294 virus database (<https://www.ncbi.nlm.nih.gov/labs/virus/>), and the polyA tail of each  
295 sequence was removed. The strains are listed in Table S2.

296

### 297 **Time-series analysis of changes in oligonucleotide compositions**

298 In the time-series analysis, the average mono- and oligonucleotide compositions (%) of  
299 viruses collected in each month were calculated for each clade. To avoid statistical  
300 fluctuations due to the small sample size, months in which fewer than 10 strains had  
301 been collected were excluded from the monthly analysis.

302

### 303 **RNA-binding motif analysis**

304 RNA-binding motifs were obtained from the ATtRACT database [30]. In this database,  
305 multiple binding motifs are registered as corresponding to one RNA-binding protein;  
306 we calculated the total number of loci containing the binding motifs for each protein in  
307 the viral genomes.

308

309

310

311

### 312 **List of abbreviations**

313 SARS-CoV-2: Severe acute respiratory syndrome coronavirus-2

314 human-CoV: human coronavirus

315 bat-CoV: bat coronavirus

316

### 317 **Ethics approval and consent to participate**

318 Not applicable

319

### 320 **Consent for publication**

321 Not applicable

322

323 **Availability of data and materials**

324 The sequence dataset analyzed in this study are stored in GISAID. Other data are  
325 available from YI.

326

327 **Competing interests**

328 The authors declared that there are no conflicts of interests.

329

330 **Funding**

331 This work was supported by JSPS KAKENHI Grant Number 18K07151, by AMED  
332 under Grant Number JP20he0622033 and by COVID-19 Counterplan Research Project  
333 (supervised by Prof. Tatsumi Hirata, NIG) from the Research Organization of  
334 Information and Systems (ROIS).

335

336 **Authors' contributions**

337 YI conceived the approach and conducted this analysis. TA developed the algorithm. TI  
338 supervised this study.

339

340 **Acknowledgements**

341 We gratefully acknowledge the authors submitting their sequences from GISAID's  
342 Database and also the valuable comments of Dr. Yashushi Hiromi of National Institute  
343 of Genetics (Mishima). We thank Springer Nature Author Services for editing this  
344 manuscript for English language.

345

346 **Figure legends**

347 **Fig. 1. Time-dependent directional changes in nucleotide compositions. (A)**

348 Average mononucleotide compositions (%) in the SARS-CoV-2 genomes of each clade  
349 isolated in each month are plotted against the elapsed month. To compare the four

350 mononucleotides, the scale widths on the vertical axis are set to the same values. The  
351 colored lines distinguishing the clade (G, GH, GR, L, V and S) are shown at the bottom  
352 of the figure. The dashed line shows the averaged compositions for all strains isolated in  
353 each month. (B) The average di- and trinucleotide compositions that primarily undergo  
354 common directional changes among the six clades are plotted against the elapsed  
355 month.

356

357 **Fig. 2. Fold changes in nucleotide composition between the epidemic start and the**  
358 **last month of analysis.** A bar plot shows the fold change in composition of each mono-  
359 or oligonucleotide; this value was calculated by dividing the nucleotide composition in  
360 the last month of analysis by that at the start of the epidemic. Each bar is colored to  
361 indicate the clade, as described in Fig. 1. Since we analyzed strains belonging to  
362 different clades separately, data from the first or last month differed among clades; see  
363 also the Methods section.

364

365 **Fig. 3. Nucleotide compositions of human and bat coronavirus sequences.** A  
366 boxplot shows the nucleotide compositions in human-CoV (alpha 229E, alpha NL63,  
367 beta HKU1 and beta OC43), bat-CoV (bat SARS, alphacoronavirus and  
368 betacoronavirus) and SARS-CoV-2 strains. Bat SARS are marked pink. A hollow arrow  
369 indicates the direction of change in oligonucleotide composition observed for SARS-  
370 CoV-2 in Figs. 1 and 2. (A) Mononucleotides. To compare the four mononucleotides,  
371 the scale widths on the vertical axis scale are set to the same values. (B) Di- and  
372 trinucleotides.

373

374 **Fig. 4. Differences in nucleotide composition between SARS-CoV-2 and human-**  
375 **CoV.** (A) Values for the square of the difference in mononucleotide composition  
376 between SARS-CoV-2 isolated in each month and human-CoV are plotted against the

377 elapsed month. The data are presented as colored or dashed lines, as described in Fig. 1.  
378 (B and C) Oligonucleotide compositions that approach and move from those of human-  
379 CoV are presented, respectively. (D) The correlation coefficients between the elapsed  
380 month from the start of the epidemic and the above differences in mono- and  
381 oligonucleotides whose directionality of change is common among six clades are  
382 presented. The results for A and G mononucleotides, which show nondirectional  
383 change, are also presented.

384

385 **Fig. 5. Time-dependent changes in the numbers of RNA-binding motif loci.** (A) The  
386 numbers of loci containing RNA-binding motifs per genome are plotted against the  
387 elapsed month. Here, we selected RNA-binding proteins for which the number of motif  
388 loci increased or decreased by at least one for all six clades from the epidemic start. The  
389 data are presented as colored or dashed lines, as described in Fig. 1A. (B) A boxplot  
390 shows the number of loci containing RNA-binding motifs in human-CoV (alpha 229E  
391 and NL63: beta HKU1 and OC43), bat-CoV (bat SARS, alphacoronavirus and  
392 betacoronavirus) and SARS-CoV-2 strains. Bat SARS are marked pink. A hollow arrow  
393 indicates the direction shown in Fig. 5A with which the oligonucleotide compositions of  
394 SARS-CoV-2 changed.

395

396

397 Table 1. Correlation coefficients for time-dependent changes in mono- and  
398 oligonucleotide compositions in SARS-CoV-2 that have increased.

399

400 Table 2. Correlation coefficients for time-dependent changes in mono- and  
401 oligonucleotide compositions in SARS-CoV-2 that have decreased.

402

403 Table 3. The number motif-containing loci for RNA-binding proteins whose  
404 occurrences have increased or decreased between strains of the first and last month of  
405 the analysis.

406

407

408 **Additional file 1**

409 Fig. S1: Average di- and trinucleotide compositions (A and B) of for SARS-CoV-2  
410 strains collected in each elapsed month.

411 Fig. S2: Oligonucleotide compositions of human and bat coronavirus sequences.

412 Fig. S3: Differences in oligonucleotide composition between SARS-CoV-2 and human-  
413 CoV.

414

415

416 **Additional file 2**

417 Table S1: List of SARS-CoV-2 strains used in the analysis.

418

419 Table S2: List of human-and bat-CoV strains used in the analysis.

420

421 Table S3: Number of SARS-CoV-2 strains in each clade isolated in each elapsed month.

422

423 Table S4: Average oligonucleotide compositions for SARS-CoV-2 strains in each clade  
424 isolated in each elapsed month.

425

426 Table S5: Correlation coefficients for time-dependent changes in oligonucleotide  
427 compositions of SARS-CoV-2.

428



429 Table S6: Fold change in compositions between strains of the first and last month of the  
430 analysis.

431 Table S7: Distance between the oligonucleotide composition of SARS-CoV-2 isolated  
432 in each elapsed month and that of human-CoV.

433

434 Table S8: Correlation coefficients for time-series changes in the distance between  
435 oligonucleotide compositions of SARS-CoV-2 and human-CoV.

436

437 Table S9: List of RNA-binding motifs.

438

439 Table S10: Numbers of motif-containing loci for RNA-binding proteins whose  
440 abundance increases or decreases between strains of the first and last month of the  
441 analysis.

442

443 Table S11: P-value from t-test to analyze the number of RNA-binding motif loci whose  
444 abundance increases or decreases between strains of the first and last month of the  
445 analysis.

446

447 Table S12: Correlation coefficients for time-dependent changes in the number of loci  
448 containing RNA-binding motifs.

449

450

## 451 **Reference**

452 1. Singhal T: A review of coronavirus disease-2019 (COVID-19). *Indian J Pediatr.*  
453 2020; 87:281-86.

454 2. García-Sastre A: Inhibition of interferon-mediated antiviral responses by influenza  
455 A viruses and other negative-strand RNA viruses. *Virology.* 2001;279: 375–84.

- 456 3. Voinnet O: Induction and suppression of RNA silencing: insights from viral  
457 infections. *Nat. Rev. Genet.* 2005;6:206–20.
- 458 4. Randall RE, Goodbourn S: Interferons and viruses: an interplay between induction,  
459 signalling, antiviral responses and virus countermeasures. *J. Gen. Virol.* 2008;89:1–  
460 47.
- 461 5. Konno Y, Kimura I, Uriu K, et al: SARS-CoV-2 ORF3b is a potent interferon  
462 antagonist whose activity is increased by a naturally occurring elongation variant.  
463 *Cell Rep.* 2020;32:108185.
- 464 6. Zhou et al: A novel bat coronavirus closely related to SARS-CoV-2 contains natural  
465 insertions at the S1/S2 cleavage site of the spike protein. *Curr Biol.* 2020;30:2196-  
466 203.
- 467 7. Nei M: *Molecular evolutionary genetics.* Columbia University Press: New York.  
468 1987.
- 469 8. Kumar S, Nei M, Dudley J, Tamura K: MEGA: a biologist-centric software for  
470 evolutionary analysis of DNA and protein sequences, *Brief Bioinform.* 2008;9:299–  
471 306.
- 472 9. Abe T, Kanaya S, Kinouchi M, et al: Informatics for unveiling hidden genome  
473 signatures, *Genome Res.* 2003;13:693–702.
- 474 10. Abe T, Sugawara H, Kinouchi M, Kanaya S, Ikemura T: Novel phylogenetic  
475 studies of genomic sequence fragments derived from uncultured microbe mixtures  
476 in environmental and clinical samples, *DNA Res.* 2005;12:281–90.
- 477 11. Iwasaki Y, Abe T, Wada K, Itoh M, Ikemura T,: Prediction of directional changes of  
478 influenza A virus genome sequences with emphasis on pandemic H1N1/09 as a  
479 model case. *DNA res* 2011;18:125-36
- 480 12. Iwasaki Y, Abe T, Wada Y, Wada K, Ikemura T: Novel bioinformatics strategies for  
481 prediction of directional sequence changes in influenza virus genomes and for  
482 surveillance of potentially hazardous strains. *BMC Infect Dis.* 2013;13:386-

- 483 13. Karlin S, Campbell AM, Mrazek J: Comparative DNA analysis across diverse  
484 genomes. *Annu. Rev. Genet.* 1998;32:185–225.
- 485 14. Wada Y, Wada K, Iwasaki Y, Kanaya S, Ikemura T: Directional and reoccurring  
486 sequence change in zoonotic RNA virus genomes visualized by time-series word  
487 count. *Sci Rep.* 2016;6:36197.
- 488 15. Wada K, Wada Y, Iwasaki Y, Ikemura T: Time-series oligonucleotide count to  
489 assign antiviral siRNAs with long utility fit in the big data era. *Gene Ther.*  
490 2017;24:668–73.
- 491 16. Wada K, Wada Y, Ikemura T: Time-series analyses of directional sequence changes  
492 in SARS-CoV-2 genomes and an efficient search method for candidates for  
493 advantageous mutations for growth in human cells. *Gene.* 2020;5:100038.
- 494 17. Qiu Y, Abe T, Nakao R, Satoh K, Sugimoto C: Viral population analysis of the  
495 taiga tick, *Ixodes persulcatus*, by using Batch Learning Self-Organizing Maps and  
496 BLAST search. *Journal of Veterinary Medical Science*, 2019;81(3):401-10.
- 497 18. Mercatelli D, Giorgi FM: Geographic and genomic distribution of SARS-CoV-2  
498 mutations. *Front Microbiol.* 2020;22:11:1800.
- 499 19. Simmonds P: Rampant C→U hypermutation in the genomes of SARS-CoV-2 and  
500 other coronaviruses: causes and consequences for their short- and long-term  
501 evolutionary trajectories. *mSphere.* 2020;24:e00408-20.
- 502 20. Paek KY, Kim CS, Park SM, Kim JH, Jang SK: RNA-binding protein hnRNP D  
503 modulates internal ribosome entry site-dependent translation of hepatitis C virus  
504 RNA. *J Virol.* 2008;82:12082-93.
- 505 21. Harris RS, Bishop KN, Sheehy AM, Craig HM, Petersen-Mahrt SK, Watt IN,  
506 Neuberger MS, Malim MH: DNA deamination mediates innate immunity to  
507 retroviral infection. *Cell.* 2003;113:803–809.

- 508 22. Mangeat B, Turelli P, Caron G, Friedli M, Perrin L, Trono D: Broad antiretroviral  
509 defence by human APOBEC3G through lethal editing of nascent reverse transcripts.  
510 Nature. 2003;424:99–103.
- 511 23. Zhang H, Yang B, Pomerantz RJ, Zhang C, Arunachalam SC, Gao L: The cytidine  
512 deaminase CEM15 induces hypermutation in newly synthesized HIV-1 DNA.  
513 Nature. 2003. 424:94–98. <https://doi.org/10.1038/nature01707>.
- 514 24. Harris RS, Dudley JP: APOBECs and virus restriction. *Virology*. 2015;479–  
515 480:131–45.
- 516 25. Sawyer SL, Emerman M, Malik HS: Ancient adaptive evolution of the primate  
517 antiviral DNA-editing enzyme APOBEC3G. *PLoS Biol*. 2004;2:E275.
- 518 26. Münk C, Willemsen A, Bravo IG: An ancient history of gene duplications, fusions  
519 and losses in the evolution of APOBEC3 mutators in mammals. *BMC Evol Biol*.  
520 2012;12:71.
- 521 27. Henry M, Terzian C, Peeters M, Wain-Hobson S, Vartanian JP: Evolution of the  
522 primate APOBEC3A cytidine deaminase gene and identification of related coding  
523 regions. *PLoS One*. 2012;7:e30036.
- 524 28. Wang W, Caldwell MC, Lin S, Furneaux H, Gorospe M: HuR regulates cyclin A  
525 and cyclin B1 mRNA stability during cell proliferation. *EMBO J*.  
526 2000;19(10):2340-50.
- 527 29. Lal A, Mazan-Mamczarz K, Kawai T, Yang X, Martindale JL, Gorospe M:  
528 Concurrent versus individual binding of HuR and AUF1 to common labile target  
529 mRNAs. *EMBO J*. 2004;23(15):3092-102.
- 530 30. Giudice G, Sánchez-Cabo F, Torroja C, Lara-Pezzi E: ATtRACT-a database of  
531 RNA-binding proteins and associated motifs. *Database (Oxford)*. 2016;7:baw035.  
532  
533  
534

535 Table 1

|     | Clade G | Clade GH | Clade GR | Clade L | Clade V | Clade S |
|-----|---------|----------|----------|---------|---------|---------|
| U   | 0.97    | 0.92     | 0.91     | 0.72    | 0.66    | 0.73    |
| UA  | 0.89    | 0.80     | 0.92     | 0.62    | 0.30    | 0.50    |
| AUU | 0.81    | 0.78     | 0.92     | 0.90    | 0.27    | 0.05    |
| CAU | 0.56    | 0.59     | 0.54     | 0.51    | 0.55    | 0.17    |
| UGU | 0.68    | 0.82     | 0.52     | 0.59    | 0.88    | 0.50    |
| UUA | 0.96    | 0.72     | 0.93     | 0.80    | 0.19    | 0.34    |
| UUG | 0.67    | 0.78     | 0.20     | 0.92    | 0.69    | 0.05    |
| UUU | 0.94    | 0.82     | 0.89     | 0.80    | 0.41    | 0.93    |

536

537 Table 2

|     | Clade G | Clade GH | Clade GR | Clade L | Clade V | Clade S |
|-----|---------|----------|----------|---------|---------|---------|
| C   | -0.95   | -0.83    | -0.95    | -0.98   | -0.97   | -0.35   |
| AG  | -0.77   | -0.71    | -0.27    | -0.57   | -0.67   | -0.44   |
| CA  | -0.73   | -0.93    | -0.85    | -0.16   | -0.45   | -0.82   |
| CC  | -0.93   | -0.87    | -0.75    | -0.90   | -0.90   | -0.09   |
| CU  | -0.81   | -0.40    | -0.79    | -0.52   | -0.15   | -0.10   |
| GA  | -0.28   | -0.90    | -0.80    | -0.62   | -0.73   | -0.29   |
| GG  | -0.60   | -0.79    | -0.65    | -0.03   | -0.10   | -0.23   |
| UC  | -0.33   | -0.23    | -0.10    | -0.58   | -0.93   | -0.41   |
| AGC | -0.57   | -0.87    | -0.41    | -0.80   | -0.78   | -0.12   |
| CCC | -0.69   | -0.81    | -0.67    | -0.94   | -0.81   | -0.50   |
| GAC | -0.73   | -0.91    | -0.65    | -0.35   | -0.18   | -0.55   |
| GAG | -0.11   | -0.58    | -0.64    | -0.50   | -0.81   | -0.02   |
| GGA | -0.73   | -0.84    | -0.75    | -0.65   | -0.81   | -0.52   |

538

539 Table 3

|         | Clade G | Clade GH | Clade GR | Clade L | Clade V | Clade S |
|---------|---------|----------|----------|---------|---------|---------|
| PTBP1   | -2.81   | -3.46    | -4.92    | -13.50  | -11.66  | 3.95    |
| HNRNPL  | -3.03   | -1.78    | -0.52    | -3.62   | -6.48   | 0.71    |
| NOVA1   | -0.04   | -2.77    | -0.70    | -3.37   | -5.04   | 1.09    |
| SRSF2   | -1.17   | -3.02    | -1.09    | -3.10   | -3.10   | 0.78    |
| ZFP36   | 1.67    | -3.50    | -1.98    | -3.48   | -4.78   | 1.86    |
| HNRNPA1 | -0.16   | -2.50    | -0.50    | -2.64   | -3.93   | 0.44    |
| ELAVL1  | 2.82    | 0.47     | 2.87     | 0.59    | 0.03    | 2.39    |
| TIA1    | -0.82   | -2.05    | -1.20    | -1.72   | -4.41   | 1.67    |
| PCBP2   | -0.37   | -2.50    | -1.03    | -2.28   | -2.35   | 0.11    |
| SRSF1   | -0.63   | -2.29    | -1.26    | -2.55   | -1.84   | 0.36    |

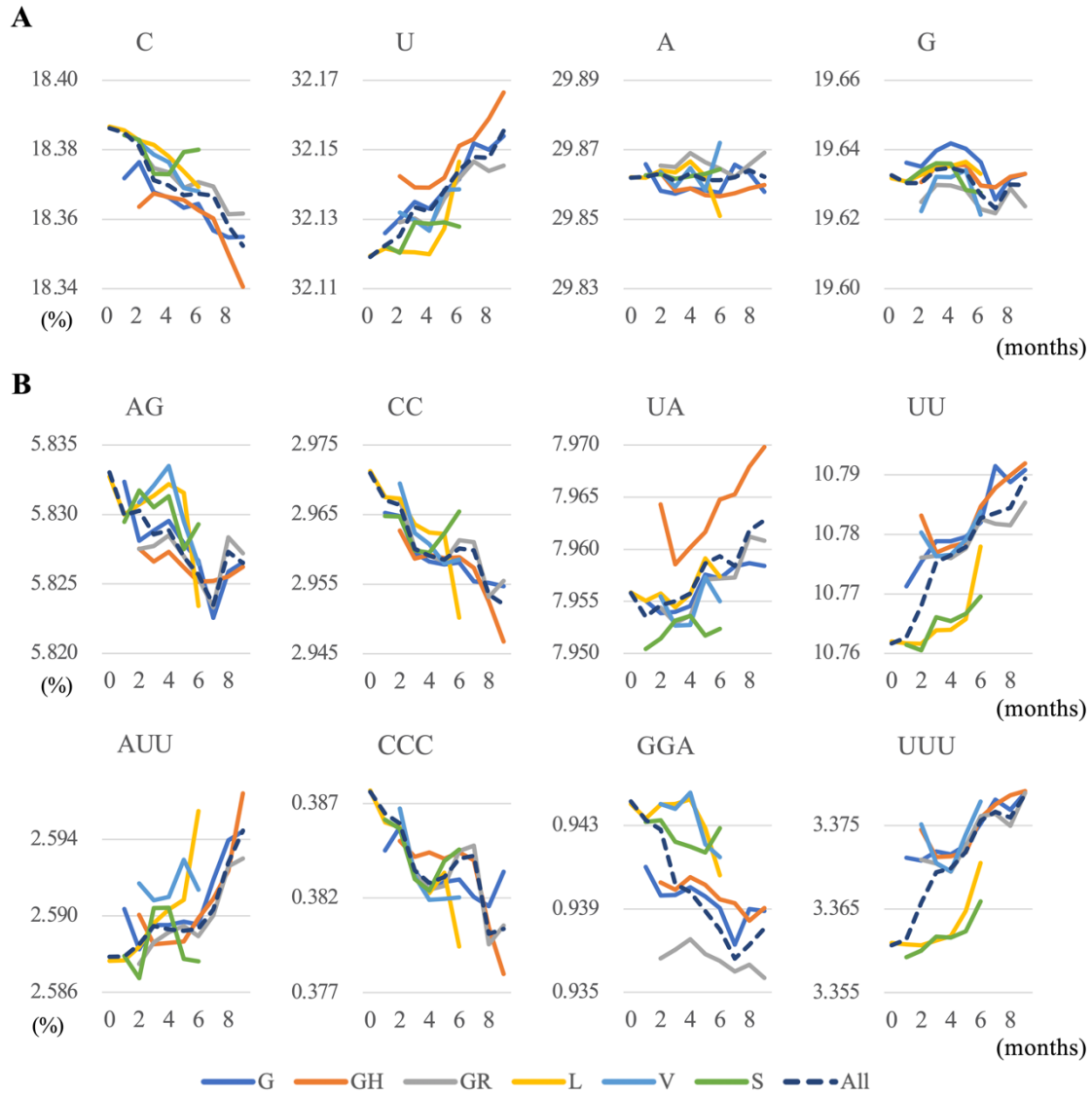
540

541

542

543

Fig. 1

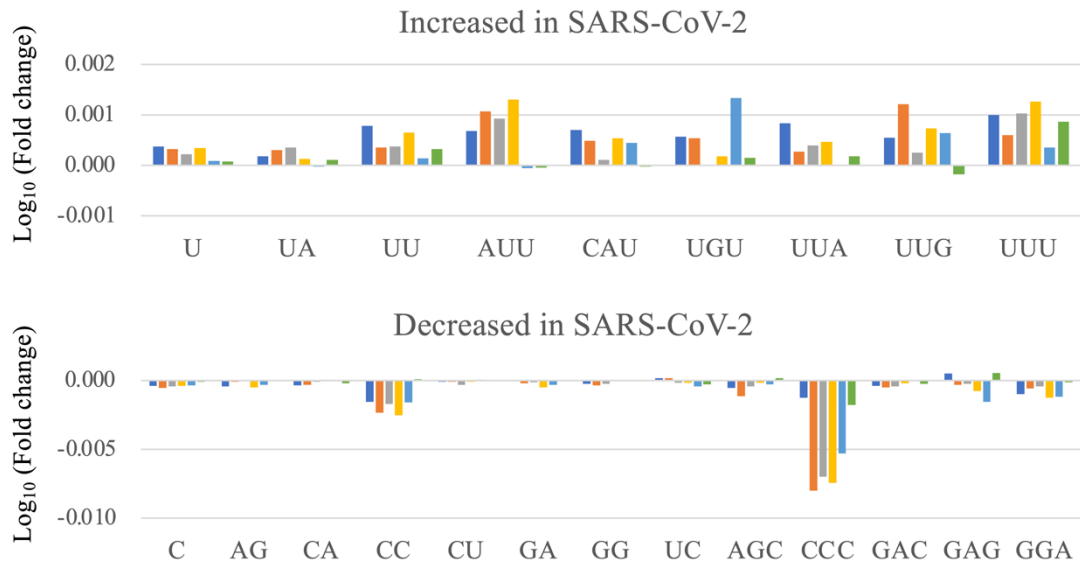


544

545

546

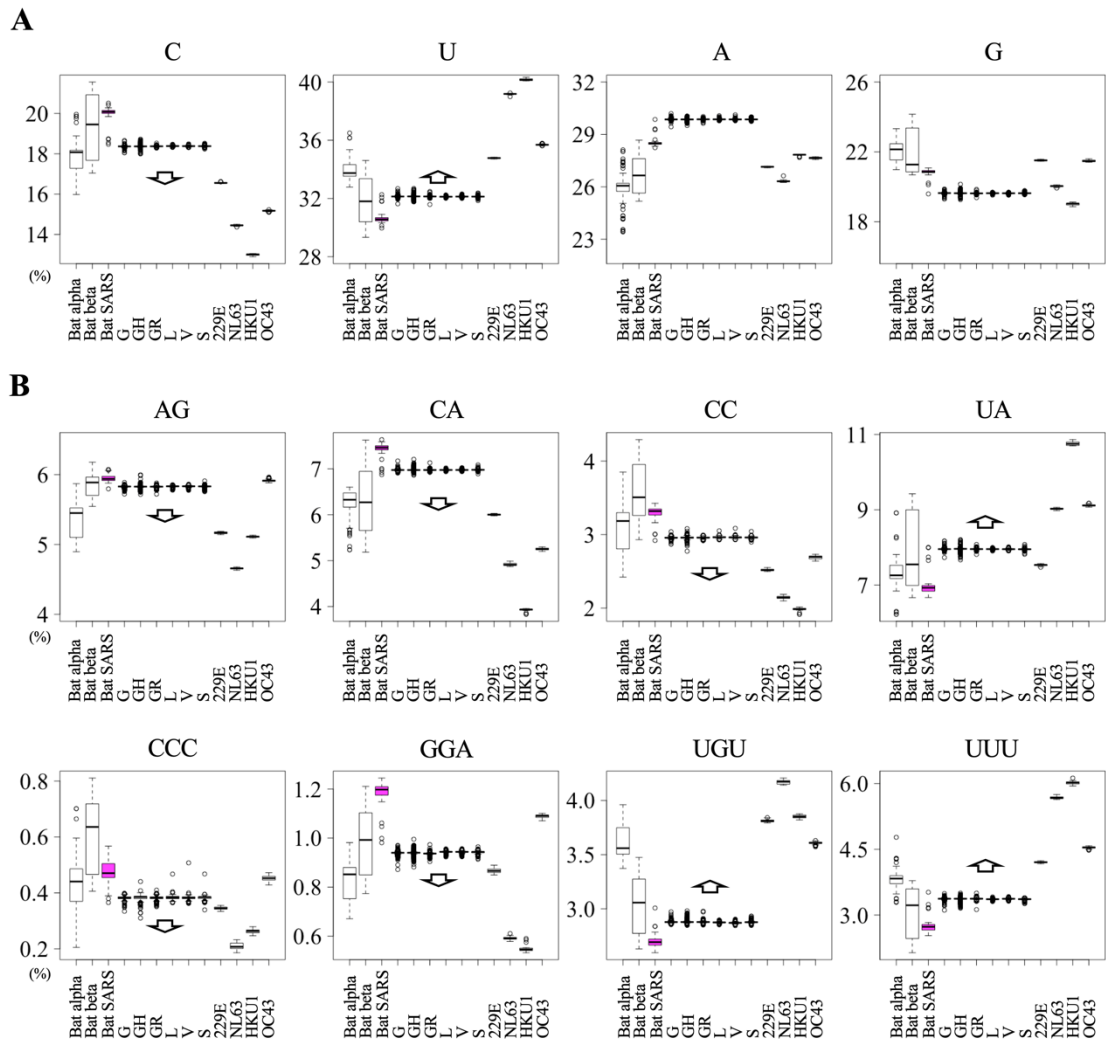
Fig. 2



547  
548



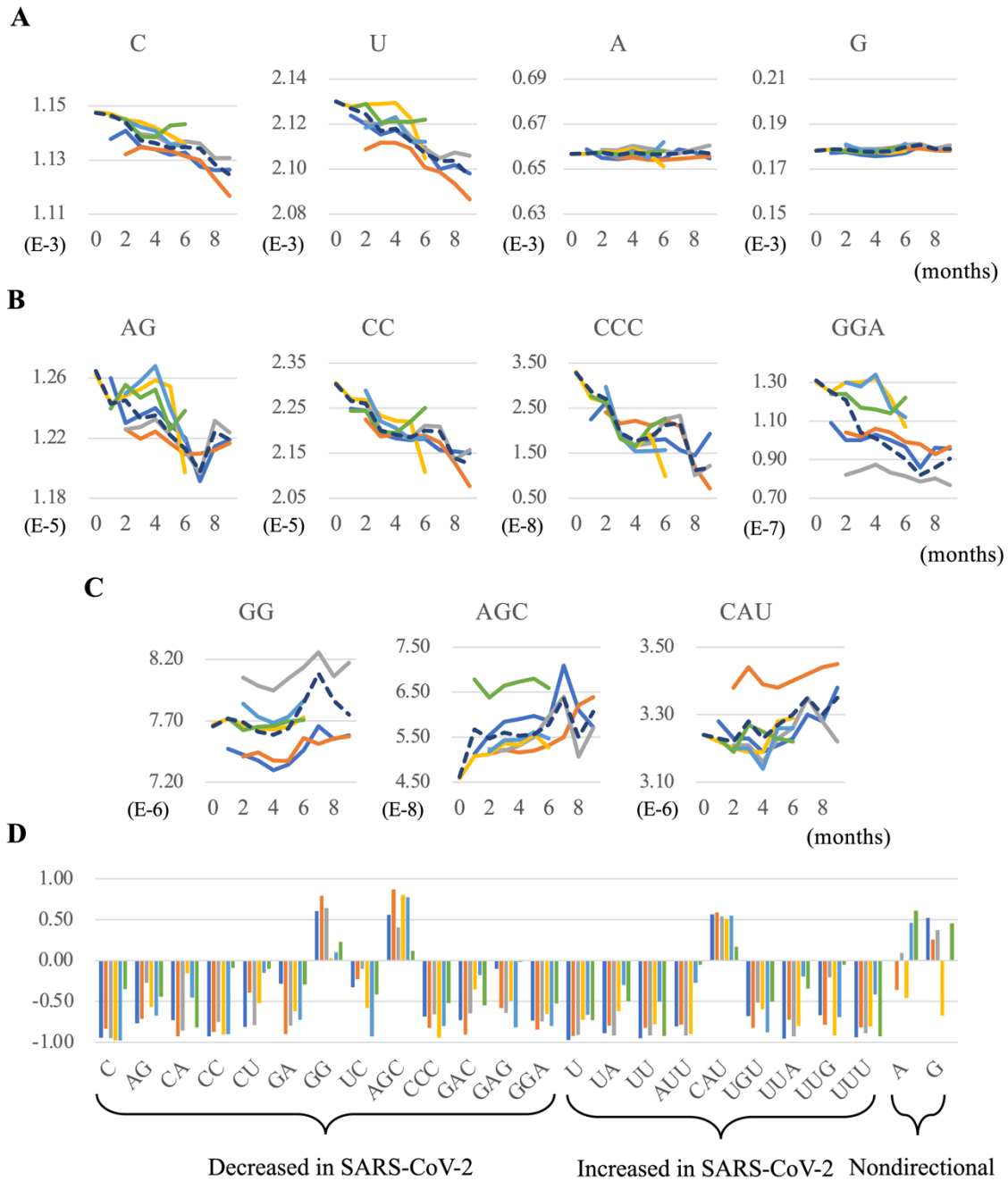
Fig. 3



549

550

Fig. 4

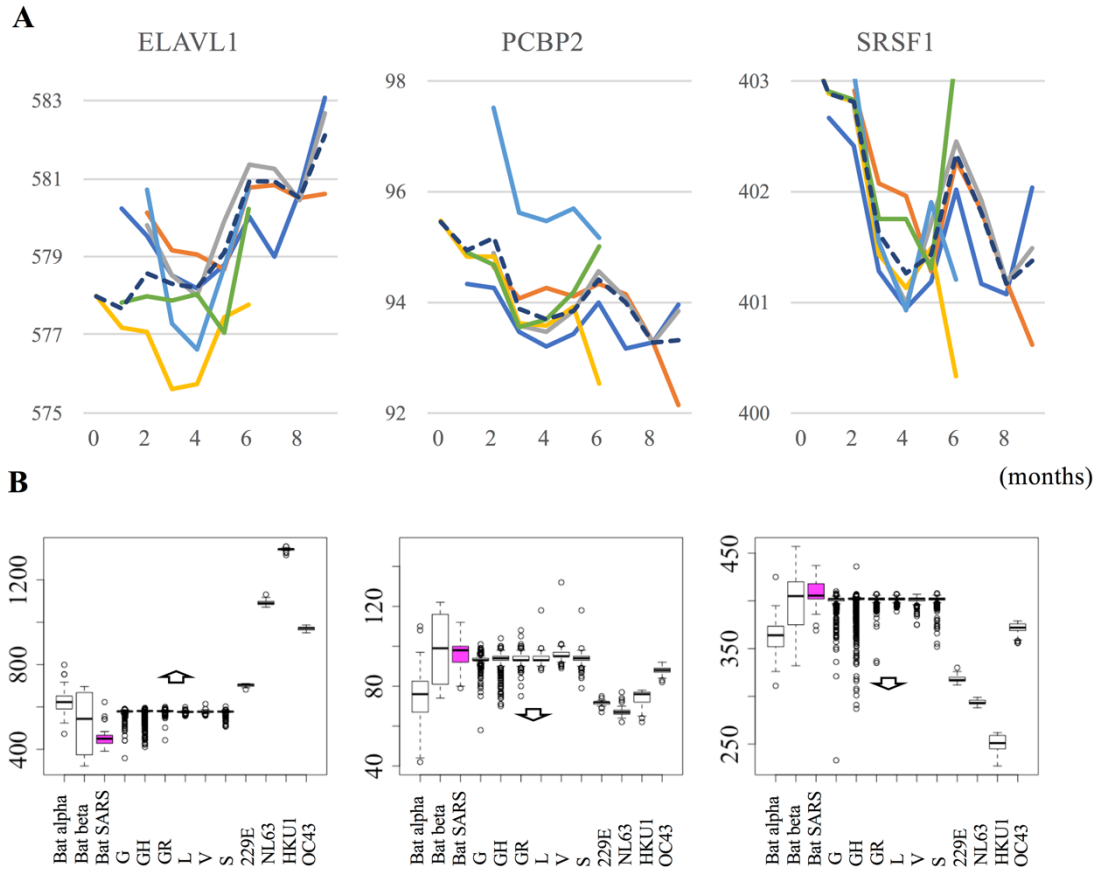


551

552

553

Fig. 5



554