

1 **Title:** Short sequence motif dynamics in the SARS-CoV-2 genome suggest a role for cytosine
2 deamination in CpG reduction.

3 **Authors:** Mukhtar Sadykov^{1*}, Tobias Mourier^{1*}, Qingtian Guan¹, Arnab Pain^{1,2,3**}

4 **Affiliation:**

5 ¹King Abdullah University of Science and Technology (KAUST), Pathogen Genomics
6 Laboratory, Biological and Environmental Science and Engineering (BESE), Thuwal-Jeddah,
7 23955-6900, Saudi Arabia;

8 ²Research Center for Zoonosis Control, Global Institution for Collaborative Research and
9 Education (GI-CoRE); Hokkaido University, N20 W10 Kita-ku, Sapporo, 001-0020 Japan;

10 ³Nuffield Division of Clinical Laboratory Sciences (NDCLS), University of Oxford, Headington,
11 Oxford, OX3 9DU, United Kingdom

12 *These authors contributed equally

13 ****Correspondence:** Arnab Pain Email: arnab.pain@kaust.edu.sa. King Abdullah University of
14 Science and Technology, Jeddah, Saudi Arabia. Phone: (+966) 54 470 0687

15 **Abbreviations:** C>U stands for cytosine to uracil substitution, the same applies to other
16 nucleotide substitutions; APOBEC - Apolipoprotein B Editing Complex; ZAP – zinc-finger
17 antiviral protein.

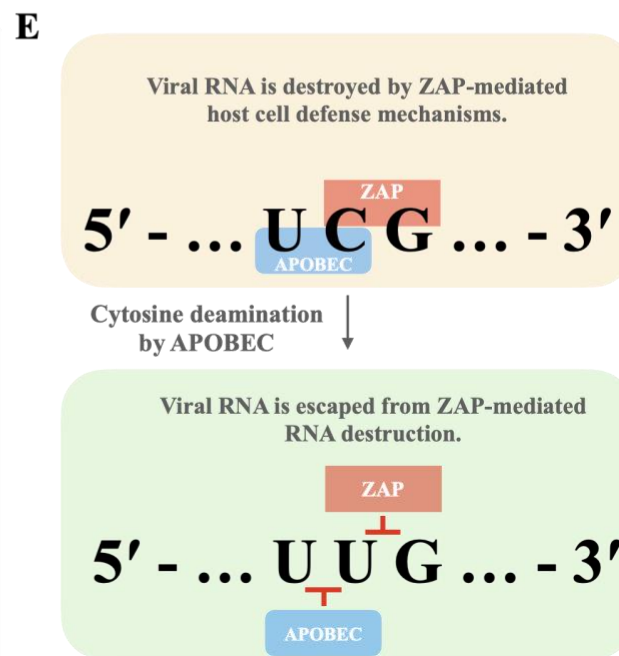
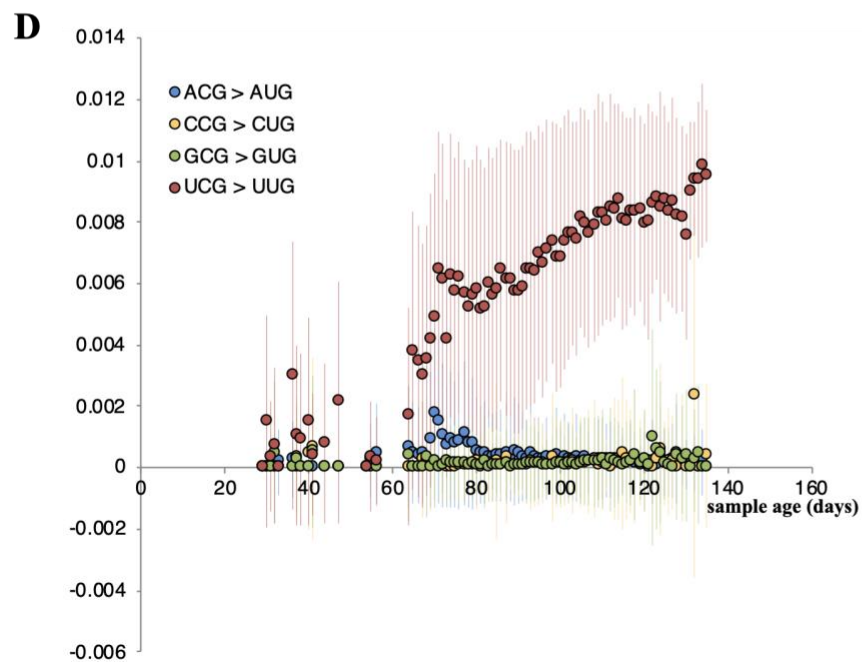
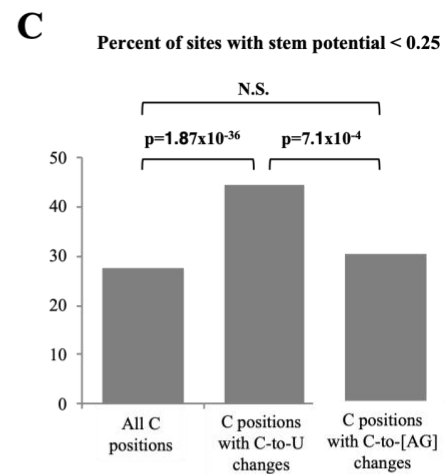
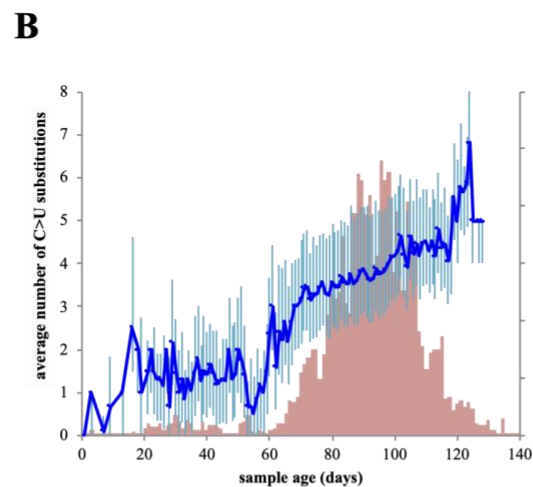
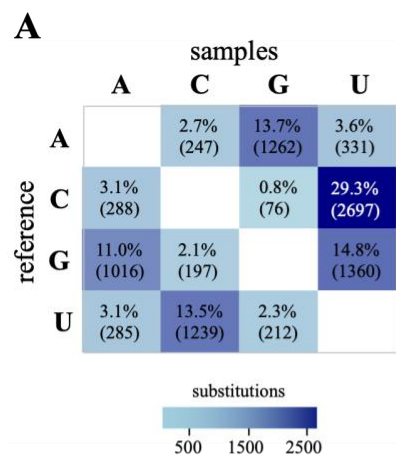
18 **Key words:** virus evolution, genome evolution, genome biology, virus-host interaction.

19 Dear Editor,

20 The APOBEC protein family are host antiviral enzymes known for catalyzing cytosine to uracil
21 deamination in foreign single-stranded DNA (ssDNA) and RNA (ssRNA) (Blanc and Davidson
22 2010; Salter and Smith 2018). Enzymatic target motifs for most of the APOBEC enzymes have
23 been experimentally identified, among which the most common were 5'-[T/U]C-3' and 5'-CC-3'
24 for DNA/RNA substrates (Salter and Smith 2018; McDaniel et al. 2020). It was recently
25 suggested that the SARS-CoV-2 undergoes genome editing by host-dependent RNA-editing
26 proteins such as APOBEC (Di Giorgio et al. 2020; Simmonds 2020; Rice et al. 2020; Schmidt et
27 al. 2020).

28

29 Given the large amount of available data and the relatively low mutation rate of the SARS-CoV-
30 2 virus (Rambaut et al. 2020), we aimed to monitor its genomic evolution on a very brief time
31 scale during the COVID-19 pandemic. Here we demonstrate progressive C>U substitutions in
32 SARS-CoV-2 genome within the timeframe of five months. We highlight the role of C>U
33 substitutions in the reduction of 5'-UCG-3' motifs and hypothesize that this progressive decrease
34 is driven by host APOBEC activity.



36 **Figure 1.** (A) SNV events observed between individual SARS-CoV-2 sample sequences
37 (n=22,164) and the reference genome. (B) The number of C>U substitutions across sample dates.
38 The average number of substitutions for each sampling day is plotted (blue line, left y-axis) with
39 plus/minus one standard deviations as error bars. The number of samples for each day is shown
40 as red bars (right y-axis). (C) Folding potential of positions with C>U changes (Supplementary
41 Text). P-values from Fisher's exact test are shown above bars. (D) The fraction of [A/C/G/U]CG
42 triplets that are changed to [A/C/G/U]UG over time. The average fractions, relative to the
43 reference genome, are shown as circles for each sampling day (x-axis). Error bars denote
44 plus/minus one standard deviation. Only dates with at least 20 samples are plotted. (E) A model
45 for the consequences of host-driven evolution by APOBEC enzymes on viral CpG dinucleotide
46 composition.

47

48 We aligned 22,164 SARS-CoV-2 genomes from GISAID to the reference genome and observed
49 a total of 9,210 single nucleotide changes with C>U being the most abundant (Figure 1A)
50 (Figure S1 & S2; Table S1; Supplementary Text). Over a period of five months, we find a steady
51 and substantial increase in C>U substitutions (Figure 1B), with almost half of them being
52 synonymous (Supplementary Text, Figure S3), and not observed for other changes (Figure S4).
53 One potential driver behind the increase in C>U changes could be the recently proposed
54 APOBEC-mediated viral RNA editing (Di Giorgio et al. 2020; Simmonds 2020) (Supplementary
55 Text). Since APOBEC3 family members display a preference for RNA in open conformation as
56 opposed to forming secondary structures (McDaniel et al. 2020), we calculated the folding
57 potential of all genomic sites that include C>U substitutions (Figure 1C). Positions with C>U
58 changes are more often located in regions with low potential for forming secondary RNA

59 structures. These observations are in agreement with the notion that members of the APOBEC
60 family are the main drivers of cytosine deamination in SARS-CoV-2 (Di Giorgio et al. 2020;
61 Simmonds 2020).

62

63 We searched for possible APOBEC genetic footprints (5'-UC-3' > 5'-UU-3') in viral dinucleotide
64 frequencies (Figure S5). Among all dinucleotides, UpC showed the highest degree of decrease,
65 while UpU exerted the highest rates of increase, which is consistent with APOBEC activity
66 (Supplementary Text).

67

68 When analyzing the context of genomic sites undergoing C>U changes, we noticed an
69 enrichment for 5'-UCG-3' motifs (Table S2). To assess the contribution of C>U changes in CpG
70 loss, we examined the dynamics of [A/C/G/U]CG trinucleotides over time (Figure 1D). The
71 progressive change (~1% over a 5-month period) of 5'-UCG-3' to 5'-UUG-3' is most striking
72 when supported by a larger number of genomes (days 70 to 115), whereas no such pattern is
73 observed for the other trinucleotides (Figure 1D). The association between cytosine deamination
74 and CpG loss is further underlined by the rapid, progressive increase in 5'-UCG-3' > 5'-UUG-3'
75 changes compared to other 5'-UC[A/C/U]-3' motifs (Figure S7). No apparent progression of 5'-
76 UCG-3' over time is observed on the negative strand, suggesting that the action of APOBEC on
77 the negative strand of SARS-CoV-2 is limited compared to the positive strand (Figure S8).

78

79 The zinc-finger antiviral protein (ZAP) selectively binds viral CpG regions that results in viral
80 RNA degradation (Takata et al. 2017). Previous studies reported that the reduced number of CpG
81 motifs in HIV and other viruses played an important role in the viral replication inside the host

82 cell, allowing the virus to escape ZAP protein activity (Takata et al. 2017). Similarly, a strong
83 suppression of CpGs is observed in SARS-CoV-2 compared to other coronaviruses (Digard et al.
84 2020). Given the high expression levels of APOBEC and ZAP genes in COVID-19 patients
85 (Blanco-Melo et al. 2020), the direct interaction of APOBEC with viral RNA (Schmidt et al.
86 2020), and our observations, we hypothesize that as a consequence of APOBEC-mediated RNA
87 editing, SARS-CoV-2 genome may escape host cell ZAP activity. Both APOBEC and ZAP are
88 interferon-induced genes that act preferentially on ssRNA in open conformation (Luo et al. 2020;
89 McDaniel et al. 2020). Initially, APOBEC and ZAP enzymes may have overlapping preferred
90 target motifs for their enzymatic functions (Figure 1E). The catalytic activity of APOBEC on 5'-
91 UC-3' leads to cytosine deamination, which destroys ZAP's specific acting site (5'-CG-3'). The
92 conversion of C>U allows viral RNA to escape from ZAP-mediated RNA destruction. Therefore,
93 uracil editing is more likely to become fixed at UCG positions due to the selective advantage this
94 conveys to subvert ZAP-mediated degradation.

95
96 Our study of sequence dynamics across the SARS-CoV-2 pandemic supplements previous
97 studies that by comparing the SARS-CoV-2 reference genome to other viral genomes address the
98 evolutionary events prior to the Wuhan SARS-CoV-2 sequence. In contrast, our approach sheds
99 light on the evolutionary events happening during the spread of SARS-CoV-2 among the human
100 population.

101
102 A recent study hypothesized that both ZAP and APOBEC provide selective pressure that drives
103 the adaptation of SARS-CoV-2 to its host (Wei et al. 2020). Here we provided one of the
104 potential mechanisms that contribute to CpG reduction in SARS-CoV-2.

105

106 In summary, our phylogeny-free approach together with other recent studies strongly support the
107 proposed model, and it merits future experimental validation. To our knowledge, this is the first
108 study linking the dynamics of viral genome mutation to two known host molecular defense
109 mechanisms, the APOBEC and ZAP proteins.

110 **Acknowledgments**

111 We thank all laboratories which have contributed sequences to the GISAID database, Zhadyra

112 Yerkes for giving her comments and helpful discussions.

113 This work was supported by funding from King Abdullah University of Science and Technology

114 (KAUST) R3T initiative. Work in AP's laboratory is supported by the KAUST faculty baseline

115 fund (BAS/1/1020-01-01).

116

117 **Author Contributions**

118 A.P. supervised the project. M.S. and T.M. designed experiments. T.M. and QG performed

119 bioinformatic analysis. M.S. wrote the draft of the manuscript. All authors discussed, edited,

120 read, and agreed to the final version of the manuscript.

121

122 **Availability of Data**

123 The data underlying this article are available in GISAID, at <https://gisaid.org>. The ID numbers of

124 genomes used are provided in Table S1.

125 **References**

- 126 Blanc V, Davidson NO. 2010. APOBEC-1-mediated RNA editing. *Wiley Interdiscip Rev Syst*
127 *Biol Med.* 2(5):594–602. doi:10.1002/wsbm.82.
- 128 Blanco-Melo D, Nilsson-Payant BE, Liu WC, Uhl S, Hoagland D, Møller R, Jordan TX, Oishi
129 K, Panis M, Sachs D, et al. 2020. Imbalanced Host Response to SARS-CoV-2 Drives
130 Development of COVID-19. *Cell.* 181(5):1036-1045.e9. doi:10.1016/j.cell.2020.04.026.
- 131 Digard P, Lee HM, Sharp C, Grey F, Gaunt E. 2020. Intra-genome variability in the dinucleotide
132 composition of SARS-CoV-2. *Virus Evol.* 6(2). doi:10.1093/ve/veaa057.
- 133 Di Giorgio S, Martignano F, Torcia MG, Mattiuz G, Conticello SG. 2020. Evidence for host-
134 dependent RNA editing in the transcriptome of SARS-CoV-2. *Sci Adv.* 6(25).
135 doi:10.1126/sciadv.abb5813.
- 136 Luo X, Wang X, Gao Y, Zhu J, Liu S, Gao G, Gao P. 2020. Molecular Mechanism of RNA
137 Recognition by Zinc-Finger Antiviral Protein. *Cell Rep.* 30(1):46-52.e4.
138 doi:10.1016/j.celrep.2019.11.116.
- 139 McDaniel YZ, Wang D, Love RP, Adolph MB, Mohammadzadeh N, Chelico L, Mansky LM.
140 2020. Deamination hotspots among APOBEC3 family members are defined by both target
141 site sequence context and ssDNA secondary structure. *Nucleic Acids Res.* 48(3):1353–1371.
142 doi:10.1093/nar/gkz1164.
- 143 Rambaut A, Holmes EC, O’Toole Á, Hill V, McCrone JT, Ruis C, du Plessis L, Pybus OG.
144 2020. A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic
145 epidemiology. *Nat Microbiol.* 5(11):1403–1407. doi:10.1038/s41564-020-0770-5.
- 146 Rice AM, Castillo Morales A, Ho AT, Mordstein C, Mühlhausen S, Watson S, Cano L, Young
147 B, Kudla G, Hurst LD. 2020. Evidence for Strong Mutation Bias toward, and Selection

- 148 against, U Content in SARS-CoV-2: Implications for Vaccine Design. *Mol Biol Evol.*
149 doi:10.1093/molbev/msaa188.
- 150 Salter JD, Smith HC. 2018. Modeling the Embrace of a Mutator: APOBEC Selection of Nucleic
151 Acid Ligands. *Trends Biochem Sci.* 43(8):606–622. doi:10.1016/j.tibs.2018.04.013.
- 152 Schmidt N, Lareau C, Keshishian H, Melanson R, Zimmer M, Kirschner L, Ade J, Werner S,
153 Caliskan N, Lander E, et al. 2020. A direct RNA-protein interaction atlas of the SARS-CoV-
154 2 RNA in infected human cells. *bioRxiv.*:2020.07.15.204404.
155 doi:10.1101/2020.07.15.204404.
- 156 Simmonds P. 2020. Rampant C→U Hypermethylation in the Genomes of SARS-CoV-2 and Other
157 Coronaviruses: Causes and Consequences for Their Short- and Long-Term Evolutionary
158 Trajectories. *mSphere.* 5(3). doi:10.1128/msphere.00408-20.
- 159 Takata MA, Gonçalves-Carneiro D, Zang TM, Soll SJ, York A, Blanco-Melo D, Bieniasz PD.
160 2017. CG dinucleotide suppression enables antiviral defence targeting non-self RNA.
161 *Nature.* 550(7674):124–127. doi:10.1038/nature24039.
- 162 Wei Y, Silke J, Aris P, Xia X. 2020. Coronavirus genomes carry the signatures of their habitats.
163 *bioRxiv.* doi:10.1101/2020.06.13.149591.