

1 **Coronavirus genomes carry the signatures of their habitats**

2 Yulong Wei^{1,+}, Jordan R. Silke^{1,+}, Parisa Aris¹, Xuhua Xia^{1,2,*}

3 1. Department of Biology, University of Ottawa, 30 Marie Curie, P.O. Box 450, Station A,
4 Ottawa, Ontario, Canada, K1N 6N5. Tel: (613) 562-5800 ext. 6886, Fax: (613) 562-5486.

5 2. Ottawa Institute of Systems Biology, Ottawa, Ontario, Canada K1H 8M5.

6 * Corresponding author E-mail: xxia@uottawa.ca

7 + Equal contribution

8

9 **ABSTRACT**

10 Coronaviruses such as SARS-CoV-2 regularly infect host tissues that express antiviral proteins
11 (AVPs) in abundance. Understanding how they evolve to adapt or evade host immune
12 responses is important in the effort to control the spread of COVID-19. Two AVPs that may
13 shape viral genomes are the zinc finger antiviral protein (ZAP) and the apolipoprotein B mRNA-
14 editing enzyme-catalytic polypeptide-like 3 protein (APOBEC3). The former binds to CpG
15 dinucleotides to facilitate the degradation of viral transcripts while the latter deaminates C into
16 U residues leading to dysfunctional transcripts. We tested the hypothesis that both APOBEC3
17 and ZAP may act as primary selective pressures that shape the genome of an infecting
18 coronavirus by considering a comprehensive number of publicly available genomes for seven
19 coronaviruses (SARS-CoV-2, SARS-CoV, MERS, Bovine CoV, Murine MHV, Porcine HEV, and
20 Canine CoV). We show that coronaviruses that regularly infect tissues with abundant AVPs have
21 CpG-deficient and U-rich genomes; whereas viruses that do not infect tissues with abundant
22 AVPs do not share these sequence hallmarks. In SARS-CoV-2, CpG is most deficient in the S
23 protein region to evaded ZAP-mediated antiviral defense during cell entry. Furthermore, over
24 four months of SARS-CoV-2 evolutionary history, we observed a marked increase in C to U
25 substitutions in the 5' UTR and ORF1ab regions. This suggests that the two regions could be
26 under constant C to U deamination by APOBEC3. The evolutionary pressures exerted by host
27 immune systems onto viral genomes may motivate novel strategies for SARS-CoV-2 vaccine
28 development.

29

30 *Running title:* Modifications in viral genomes by host

31 *Key words:* SARS-CoV-2; APOBEC3; ZAP; viral evolution

32

33

34 INTRODUCTION

35 The emergence of SARS-CoV-2 pandemic poses a serious global health emergency.
36 Understanding how coronaviruses adapt or evade tissue-specific host immune responses is,
37 therefore, important in the effort to control the spread of COVID-19 and to facilitate vaccine-
38 development strategies. As obligate parasites, coronaviruses evolve in mammalian hosts and
39 carry genomic signatures shaped by their host-specific environments. At the tissue level,
40 mammalian species provide different cellular environments with varying levels of antiviral and
41 RNA modification activity. Two antiviral proteins (AVPs) that may contribute to the modification
42 of viral genomes are the zinc finger antiviral protein (ZAP, gene name ZC3HAV1 in mammals)
43 and the apolipoprotein B mRNA-editing enzyme-catalytic polypeptide-like 3 (APOBEC3), both of
44 which exhibit tissue-specific expression (FAGERBERG *et al.* 2014).

45 ZAP is a key component in the mammalian interferon-mediated immune response that
46 specifically targets CpG dinucleotides in viral RNA genomes (MEAGHER *et al.* 2019) to inhibit viral
47 replication and signal for viral genome degradation (FICARELLI *et al.* 2020; GUO *et al.* 2007;
48 MEAGHER *et al.* 2019; TAKATA *et al.* 2017). This host immune response acts against not only
49 retroviruses such as HIV-1 (FICARELLI *et al.* 2020; ZHU *et al.* 2011), but also single-stranded RNA
50 viruses such as Ecovirus 7 (ODON *et al.* 2019), Zika virus (TRUS *et al.* 2020), and Influenza virus
51 (GREENBAUM *et al.* 2008). It follows that cytoplasmic ZAP activity should impose a strong
52 avoidance of CpG dinucleotides in RNA viruses that target host tissues abundant in ZAP. For
53 instance, while HIV-1 infects lymph organs where ZAP is abundant (FAGERBERG *et al.* 2014), its
54 genome is also strongly CpG-deficient. Notably, the viral fitness of HIV-1 has been shown to
55 diminish as its genomic CpG content increases within a sample of patients (THEYS *et al.* 2018).
56 Many other pathogenic single-stranded RNA viruses, including coronaviruses, also exhibit
57 strong CpG deficiency (ATKINSON *et al.* 2014; GREENBAUM *et al.* 2008; GREENBAUM *et al.* 2009;
58 TAKATA *et al.* 2017; YAP *et al.* 2003), but selection for CpG deficiency disappears in ZAP-deficient
59 cells (TAKATA *et al.* 2017).

60 The ZAP-mediated RNA degradation is cumulative (TAKATA *et al.* 2017). When CpG dinucleotides
61 are added to individual viral segment 1 or 2 in HIV-1, the inhibitory effect of ZAP is weak.
62 However, when the same CpG dinucleotides are added to both segments 1 and 2, the ZAP
63 inhibition effect is strong (TAKATA *et al.* 2017). This implies that only RNA sequences of sufficient
64 length would be targeted by ZAP. Thus, although SARS-CoV-2 has the lowest genomic CpG
65 contents found in Betacoronaviruses (XIA 2020), only ORF1ab and spike (S) mRNAs are likely
66 targets by ZAP. It is therefore not surprising that these mRNAs, especially the S mRNA, exhibit
67 lower CpG than other shorter genes (DI GIOACCHINO *et al.* 2020; KIM *et al.* 2020).

68 Aside from ZAP, the APOBEC3 cytidine deaminase enzymes have garnered substantial attention
69 for their role in the antiviral immune response (CULLEN 2006; HARRIS and DUDLEY 2015). Through

70 a mechanism largely derived from HIV-1 studies, APOBEC3 enzymes have been prominently
71 reported to disrupt the structure and function of HIV-1 viruses by hypermutating minus strand
72 viral cDNA, causing defects in the viral transcript and inhibiting reverse transcription (CHIU and
73 GREENE 2008; HARRIS and DUDLEY 2015; HAYWARD *et al.* 2018; NABEL *et al.* 2013; RODRIGUEZ-FRIAS *et al.*
74 *et al.* 2013). For instance, APOBEC3G (HARRIS *et al.* 2003; MANGEAT *et al.* 2003) and 3F (ZHENG *et al.*
75 2004) catalyzes C to U deamination at the HIV-1 minus-strand DNA during reverse transcription,
76 this triggers G to A hypermutation in the nascent retroviral DNA. HIV-1 avoids these deleterious
77 effects by expressing Vif, a protein which targets and degrades APOBEC3 enzymes (SHEEHY *et al.*
78 2002; WANG and WANG 2009). Despite these findings, the possibility that APOBEC3 paralogues
79 may act directly to edit ssRNA viruses has not been widely explored but the potential should
80 not be excluded, especially for viruses that do not encode a Vif analogue such as SARS-CoV-2.

81 Indeed, APOBEC3 is now known to modify a variety of RNA sequences. For instance, RNA
82 binding activity facilitates the packaging of APOBEC3 into virions (BOGERD and CULLEN 2008;
83 ZHANG *et al.* 2010). Furthermore, C to U RNA editing has been demonstrated in macrophages,
84 monocytes, and lymphocytes by both APOBEC3A and 3G (SHARMA *et al.* 2016; SHARMA *et al.*
85 2015; SHARMA *et al.* 2019). Additionally, APOBEC3C, 3F, and 3H may inhibit HCoV-NL63
86 coronavirus infection in humans, yet it remains unclear whether C to U RNA editing was
87 involved (MILEWSKA *et al.* 2018). More recently, studies have shown evidence that SARS-CoV-2
88 genomes are driven towards increasing U content and decreasing C content (DI GIORGIO *et al.*
89 2020; JIANG 2020; SIMMONDS 2020; VICTOROVICH *et al.* 2020). Resultantly, the possibility of C-U
90 editing by APOBEC3 at the RNA level could effectively disrupt the structure and protein function
91 of positive single-stranded RNA viruses.

92 We hypothesize that both APOBEC3 and ZAP act in concert as the primary selective pressure
93 driving the adaptation of an infecting coronavirus over the course of its evolutionary history in
94 specific host tissues. To test this hypothesis, we examined which antiviral proteins are effective
95 against coronaviruses and how the immune response is subverted. We predict that when a
96 virus regularly infects host tissues that are deficient in AVPs, there will be no strong directional
97 substitutions resulting in decreased CpG dinucleotides or elevated U residues, as these
98 evolutionary forces will be weak when ZAP and APOBEC3 are lowly expressed. Conversely,
99 when a virus regularly infects host tissues that are abundant in AVPs, these antiviral responses
100 will exert their influence on viral genomes. Consequently, viral genomes should tend towards
101 reduced CpG dinucleotides to elude ZAP-mediated cellular antiviral defense, and increased U
102 residues because of RNA editing by APOBEC3 proteins.

103 Our investigation considers a comprehensive number of publicly available genomes for seven
104 coronaviruses (the Betacoronaviruses SARS-CoV-2, SARS-CoV, MERS, Bovine CoV, Murine MHV,
105 and Porcine HEV, and the Alphacoronavirus Canine CoV,) as well as studies with tissue-specific

106 ZAP and APOBEC3 gene expressions in five host species (human, cattle, dog, mice, and pig). We
107 found that all surveyed coronaviruses regularly infect tissues with high mRNA expressions of
108 both ZAP and APOBEC3, except Murine MHV. Expectedly, all surveyed coronaviruses, except
109 Murine MHV, have high global CpG deficiency, with SARS-CoV-2 genomes having the lowest
110 CpG content. More specifically, we observed a nonuniform distribution of CpG content across
111 12 SARS-CoV-2 viral regions and noted that CpG is most deficient in the region encoding the S
112 protein that mediates cell entry by ACE2 binding. Taken together, these observations suggest
113 SARS-CoV-2 has evolved in a tissue with high ZAP expression, and its persistence indicates that
114 it has successfully evaded the ZAP-mediated antiviral defense during cell entry.

115 In line with evidence of RNA-level C to U deamination by APOBEC3 enzymes (BISHOP *et al.* 2004;
116 SHARMA *et al.* 2016; SHARMA *et al.* 2015; SHARMA *et al.* 2019), Bovine CoV, Canine CoV, and
117 Porcine HEV all exhibit high global U content and low global C content whereas the genomes of
118 Murine MHV and the much more recent human coronaviruses (SARS-CoV-2, SARS-CoV, and
119 MERS) exhibit notably lower U content and higher C content. To elucidate the early stages of
120 SARS-CoV-2 genomic evolution, we analyzed the patterns of single nucleotide polymorphisms
121 (SNPs) in local viral regions of complete genomes that were collected in the span of four
122 months (from December 31, 2019 to May 6, 2020) since SARS-CoV-2 was initially isolated. We
123 observed that the occurrence of C to U substitutions is strikingly more prevalent than any other
124 SNPs, especially in the 5' UTR and ORF1ab regions. This suggests that the 5' UTR and ORF1ab
125 regions are under constraint by these enzymes. Indeed, both APOBEC3 and ZAP exert selective
126 pressure on the RNA genome compositions of coronaviruses that regularly infect tissues
127 expressing the two antiviral genes in abundance, but they do not affect the RNA genomes of
128 viruses that avoid infecting tissues with high antiviral gene expression.

129

130 **MATERIALS AND METHODS**

131 **Retrieving and processing the *APOBEC3* and *ZAP* genes and their tissue specific gene** 132 **expressions in five mammalian species**

133 The NCBI Nucleotide Database was queried for all records containing “APOBEC3” and
134 “ZC3HAV1L” as gene names, “Mammalia” as a taxonomic class, and “*Homo sapiens*”, “*Bos*
135 *taurus*”, “*Canis lupus familiaris*”, “*Mus musculus*”, and “*Sus scrofa*” as species. These five
136 species were selected because they have extensive tissue-specific gene expression studies (as
137 discussed below). Next, each entry was searched for /product= ‘apolipoprotein B mRNA editing
138 enzyme, catalytic polypeptide 3’ and for /product= ‘zinc finger CCCH-type containing, antiviral
139 1’, whole-genome and chromosome-wide results were excluded, and only the coding DNA

140 sequence region of APOBEC3 and ZC3HAV1 isoforms were extracted in FASTA format along with
141 their ENSEMBL Accession IDs.

142 To compare gene expressions of APOBEC3 and ZC3HAV1L among tissues, we retrieved publicly
143 available RNA Sequencing and Microarray studies that each sampled at least 10 mammalian
144 tissues. The five mammalian species that have extensive tissue-specific mRNA expressions are
145 *Homo sapiens*, *Bos taurus*, *Canis lupus familiaris*, *Mus musculus*, and *Sus scrofa*. For *Homo*
146 *sapiens*, tissue-specific mRNA expressions were retrieved in averaged FPKM values from all 171
147 RNA-Seq datasets in BioProject PRJEB4337(FAGERBERG *et al.* 2014), 48 RNA-Seq datasets in
148 BioProject PRJEB2445, 20 RNA-Seq datasets in BioProject PRJNA280600 (DUFF *et al.* 2015), and
149 in median TPM values from all RNA-Seq datasets available in the GTEx Portal (LONSDALE *et al.*
150 2013). For *Mus musculus*, tissue-specific mRNA expressions were retrieved in averaged FPKM
151 values from all 741 RNA-Seq datasets in BioProject PRJNA66167 (mouse ENCODE consortium)
152 (YUE *et al.* 2014) and in average TPM values from all 79 RNA-Seq datasets in BioProject
153 PRJNA516470 (NAQVI *et al.* 2019). For *Sus scrofa*, tissue-specific mRNA expressions were
154 retrieved in averaged FPKM values from TISSUE 2.0 integrated datasets (PALASCA *et al.* 2018).
155 For *Canis lupus familiaris*, tissue-specific gene expressions were retrieved in averaged
156 fluorescence intensity units (FIU) from all 39 microarray datasets in BioProject PRJNA124245
157 (BRIGGS *et al.* 2011), and in averaged TPM values from all 75 RNA-Seq datasets in BioProject
158 PRJNA516470 (NAQVI *et al.* 2019). Lastly, for *Bos taurus*, tissue-specific mRNA expressions were
159 retrieved in averaged FPKM values from 42 RNA-Seq datasets in the Bovine Genome Database
160 (SHAMIMUZZAMAN *et al.* 2019).

161 Given that the data extracted were from multiple independent sources, thus not directly
162 comparable, the relative mRNA expression level designations (high or low) for APOBEC3 and
163 ZAP isoforms in a given tissue were derived from comparisons among AVP expressions in all
164 tissues in each independent source. Specifically, we calculated the proportion of mRNA
165 expression (PME) as:

$$166 \quad PME = \frac{\text{mRNA expression value in a specific tissue}}{\text{summed mRNA expression values in all tissues}} \quad (1)$$

167 PME values were calculated from averaged TPM values in 24 human tissues using all RNA-Seq
168 datasets available in the GTEx Portal (LONSDALE *et al.* 2013), from averaged FPKM values in 26
169 cattle tissues using the Bovine Genome Database (SHAMIMUZZAMAN *et al.* 2019), from averaged
170 FPKM values in 33 pig tissues using TISSUE 2.0 integrated datasets (PALASCA *et al.* 2018), from
171 averaged FPKM values in 17 mice tissues using all 741 RNA-Seq datasets in mouse ENCODE
172 consortium (YUE *et al.* 2014), from averaged FPKM values in 12 mice tissues using 79 RNA-Seq
173 datasets in BioProject PRJNA516470 (NAQVI *et al.* 2019), and from averaged fluorescence
174 intensity units in 10 dog tissues using all 39 microarray datasets in BioProject PRJNA124245

175 (BRIGGS *et al.* 2011). Next, we calculated the averaged PME value by considering all tissue-
176 specific PME values in each independent source. Finally, for each AVP, tissue-specific PMEs
177 were designated as high if they are greater than the averaged PME value and low if they are
178 less than the averaged PME. In addition, each column in Supplemental figures S1 and S2 with
179 column title designations “APOBEC3” or “ZC3HAV1” contains the tissue-specific AVP
180 expressions from an individual source, where darkest blue represents the tissue with the
181 highest mRNA expression and darkest red represents the lowest mRNA expression.

182 **Retrieving and processing the genomes and regular habitats of coronaviruses infecting five** 183 **mammalian species**

184 The genome, Accession ID, and Sample Collection Date of 28475 SARS-CoV-2 samples were
185 retrieved from the China National Center for Bioinformation (CNCB)
186 (<https://bigd.big.ac.cn/ncov/variation/statistics?lang=en>, last accessed May 16, 2020), among
187 which 2666 strains were selected because they were annotated as having complete genome
188 sequences and high sequencing quality. Additionally, the complete genomic sequences of 403
189 MERS strains, 134 SARS-CoV strains, 20 Bovine CoV strains, 2 Canine CoV strains, 26 Murine
190 HEV strains, and 10 Porcine HEV strains were downloaded from the National Center for
191 Biotechnology Information (NCBI) Nucleotide Database (<https://www.ncbi.nlm.nih.gov/>).

192 We computed the nucleotide and di-nucleotide frequencies in each viral genome. Among
193 strains of the same coronavirus, some genomic sequences have long poly-A tails that are
194 missing in other sequences. Some also have a longer 5' untranslated region (5' UTR) than
195 others. To make a fair comparison between strains, the genomes were aligned with MAFFT
196 version 7 (KATO and STANDLEY 2013), with the slow but accurate G-INS-1 option for 134 SARS-
197 CoV, 20 Bovine CoV, 2 Canine CoV, 26 Murine MHV, and 10 Porcine HEV strains, and with the
198 fast FFT-NS-2 option for large alignments for 2666 SARS-CoV-2 and 403 MERS strains. Next,
199 using the ‘Sequence Manipulation’ tool in DAMBE7 (XIA 2018), the 5' UTR sequences were
200 trimmed away until the first fully conserved nucleotide position. Similarly, the 3' UTR sequences
201 were trimmed out up to the last fully conserved nucleotide position. Then, gaps were removed
202 from each trimmed genome, and the global nucleotide and dinucleotide frequencies as well as
203 their respective proportions (denoted as $P_{\text{Nucleotide}}$ and $P_{\text{Dinucleotide}}$, respectively) were computed
204 in DAMBE under “Seq. Analysis|Nucleotide & di-nuc Frequency”. Additionally, nucleotide and
205 di-nucleotide frequencies were similarly computed for whole, untrimmed, genomes. Finally, the
206 conventional index of CpG deficiency (CARDON *et al.* 1994; KARLIN *et al.* 1997) was calculated as:

$$207 \quad I_{\text{CpG}} = \frac{P_{\text{CG}}}{P_{\text{C}}P_{\text{G}}} \quad (2)$$

208 The index is expected to be 1 with no CpG deficiency or excess, smaller than 1 if CpG is deficient
209 and greater than 1 if CpG is in excess.

210 Next, among 2666 high sequence quality and complete SARS-CoV-2 genomes from CNCB, we
211 randomly selected one genome from each collection date, inclusively between December 31,
212 2019 (first isolate) and May 6, 2020 (most recent isolate, database last accessed on May 16,
213 2020), that have complete records of local region annotations and nucleotide sequences in
214 NCBI. A total of 99 variants (or samples) were retrieved across 127 days since SARS-CoV-2
215 (strain Wuhan-Hu-1, MN908947) was first sequenced. For each of the 99 samples, the
216 nucleotide sequence of 12 out of 13 viral regions (5' UTR, ORF1ab, S, ORF3, E, M, ORF6, ORF7a,
217 ORF8, N, ORF10, and 3' UTR) were extracted from DAMBE in FASTA format, and local nucleotide
218 and dinucleotide frequencies and their proportions were computed for each region. ORF7b was
219 omitted from the analysis because it was not annotated in 30 out of 99 samples, including the
220 reference genome Wuhan-Hu-1 (MN908947). To determine the sequence mutation patterns
221 over time at each viral region, the nucleotide sequences from our 99 genomes were first
222 aligned with MAFFT with the slow but accurate G-INS-1 option; then each aligned sequence was
223 pair-wise assessed for single nucleotide polymorphisms (SNPs) using DAMBE's "Seq.
224 Analysis|Nucleotide substitution pattern" with reference genome = Wuhan-Hu-1 (MN908947,
225 first sequenced) and Default genetic distance = F84.

226 Host tissues that are infected by SARS-CoV-2, MERS, and SARS-CoV in human, Bovine CoV in
227 cattle, Canine CoV in dog, Murine HEV in mice, and Porcine HEV in pig were identified through
228 an exhaustive large-scale manual search on relevant evidence-based primary source studies.
229 The studies considered included clinical course, autopsy, and experimental infections, but
230 cross-host studies were excluded. In total, tissue infections were determined from 25 SARS-CoV
231 studies, 11 SARS-CoV-2 studies, 8 MERS studies, 15 Murine CoV (MHV) studies, 9 Porcine HEV
232 studies, 18 Canine CoV studies, and 10 Bovine CoV studies (Supplemental File S1). Resultantly,
233 the regular tissue habitats of viruses were determined based on the prevalence of virus
234 detection in host tissues across studies. For example, among studies on SARS-CoV-2, some
235 tissue infections (e.g., infections in the lung and intestine) are frequently recorded while others
236 are rarely recorded (e.g., stomach). To compare the relative prevalence of SARS-CoV-2
237 infection, in the lung vs. in other tissues for instance, we calculated commonness of detection
238 (COD) as:

$$239 \quad COD = \frac{\text{number of times lung infection is recorded}}{\text{number of recorded infections in all tissues}} \quad (3)$$

240 DATA AVAILABILITY

241 Supplemental File S1 contains reference compilation of virus regular habitats. Supplemental File
242 S2 and S3 contain nucleotide and di-nucleotide frequencies in trimmed and whole viral
243 genomes, respectively. Supplemental File S4 contains the global and local mutation patterns in
244 SARS-CoV-2 genomes. Supplemental File S5 contains the local CpG dinucleotide frequencies in a

245 sample of 99 SARS-CoV-2 genomes. Supplemental File S6 contains Supplemental figures S1 to
246 S5.

247

248 **RESULTS**

249 **All host-specific coronaviruses except Murine MHV regularly infect host tissues that highly** 250 **express AVPs**

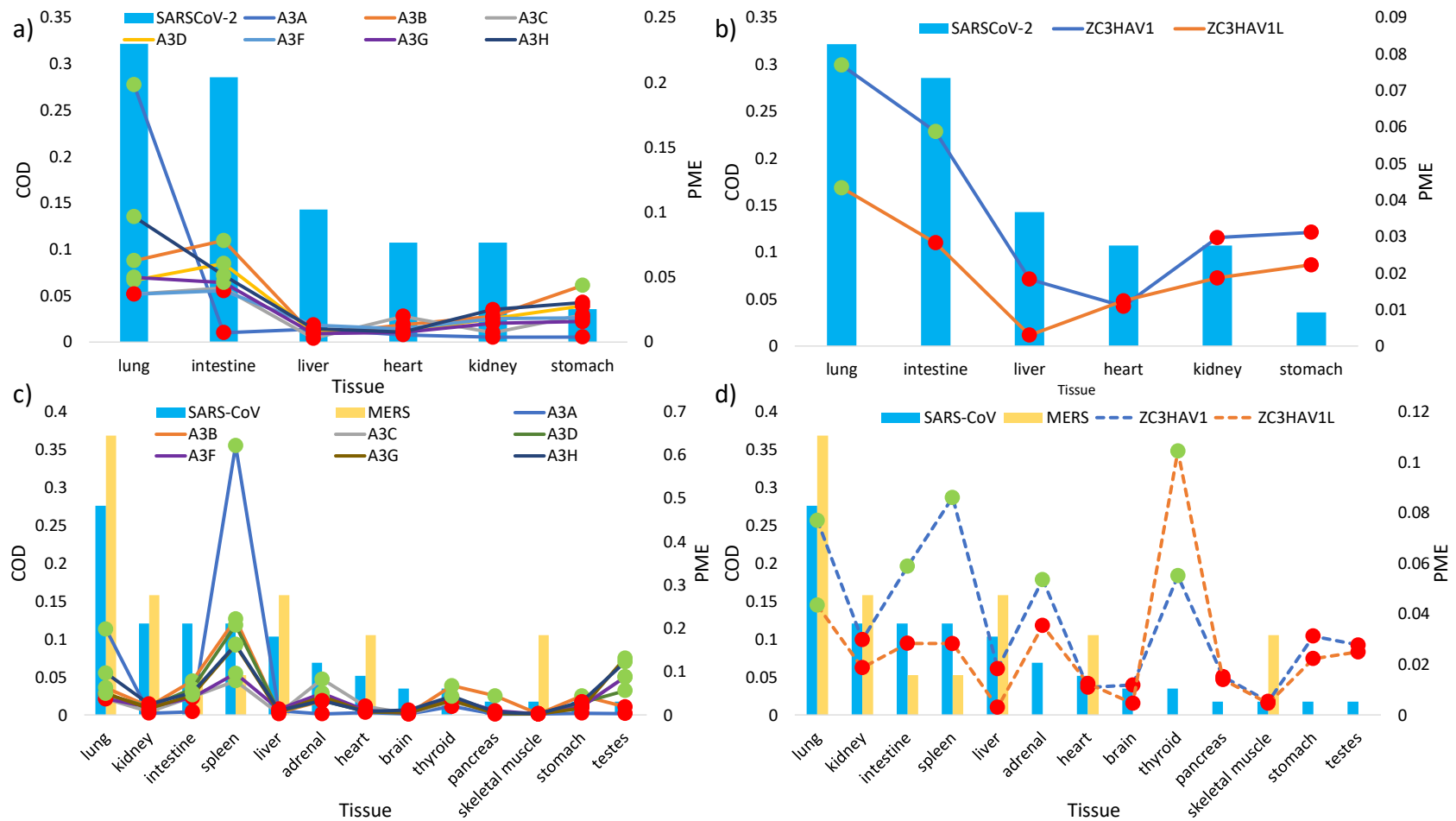
251 The tissue-specific mRNA expressions of 7 human APOBEC3 gene isoforms (A3A, A3B, A3C, A3D,
252 A3F, A3G, and A3H) and 2 human ZAP isoforms (ZC3HAV1 and ZC3HAV1L) were retrieved from
253 publicly available RNA-Seq datasets (see Materials and Methods) and averaged FPKM values
254 were compared within study. Supplemental Fig. S1 shows that all 3 human coronaviruses infect
255 the lung, heart, liver, kidney, and stomach; SARS-CoV and SARS-CoV-2 additionally infect the
256 intestines; SARS-CoV and MERS additionally infect the skeletal muscle; but only SARS-CoV has
257 been reported to infect the lymph node and the spleen (Supplemental Fig. S1; references with
258 records of tissue infection are located in Supplemental File S1).

259 We determined which human tissues are commonly infected by coronaviruses and whether
260 these tissues express AVPs in abundance. Figure 1 shows the commonness of detection of
261 SARS-COV-2, SARS-CoV, and MERS (CODs) (see Equation 3, Materials and Methods) in human
262 tissues and the relative proportions of mRNA expression (PME) of APOBEC3 and ZAP isoforms
263 (see Equation 1, Materials and Methods) in each tissue that is susceptible to infection.
264 Furthermore, in each susceptible tissue, the relative mRNA expressions of AVPs (in PME values)
265 were determined as high (in green) or low (in red) (see Materials and Methods). Out of the 6
266 tissues with records of infection, the lung and the intestine are regular habitats of SARS-CoV-2
267 with the highest CODs. Furthermore, both tissues contain high PMEs for some APOBEC3
268 isoforms (Fig. 1a: A3A, A3B, A3D, A3G, A3H in the lung, and A3B, A3D, A3G, and A3H in the
269 intestine) and for ZC3HAV1 (Fig. 1b). In the 4 tissues (liver, heart, kidney, and stomach) where
270 infection is less frequently observed (less COD values), APOBEC3 and ZAP PMEs are also low
271 (Fig. 1a, 1b). Similarly, some regular habitats of SARS-CoV-2 (lung, kidney, intestine, spleen,
272 liver) and of MERS (lung, kidney, and liver) also host high PMEs for some APOBEC3 isoforms
273 (Fig. 1c) and for some ZAP isoforms(Fig. 1d). Hence, all 3 human coronaviruses can regularly
274 infect host tissues that express relatively high AVP expressions and display no strong preference
275 for tissues with AVP deficiency.

276 Similarly, we retrieved averaged mRNA expression levels of AVP isoforms in four other
277 mammalian species (cattle, dog, pig, mice) and reference records of tissue-specific infections of
278 their coronaviruses (Supplemental File S1, Supplemental Fig. S2). We determined the regular
279 habitats (by COD) for Porcine HEV in pig (Fig. 2a), Canine CoV in dog (Fig. 2b), Bovine CoV in

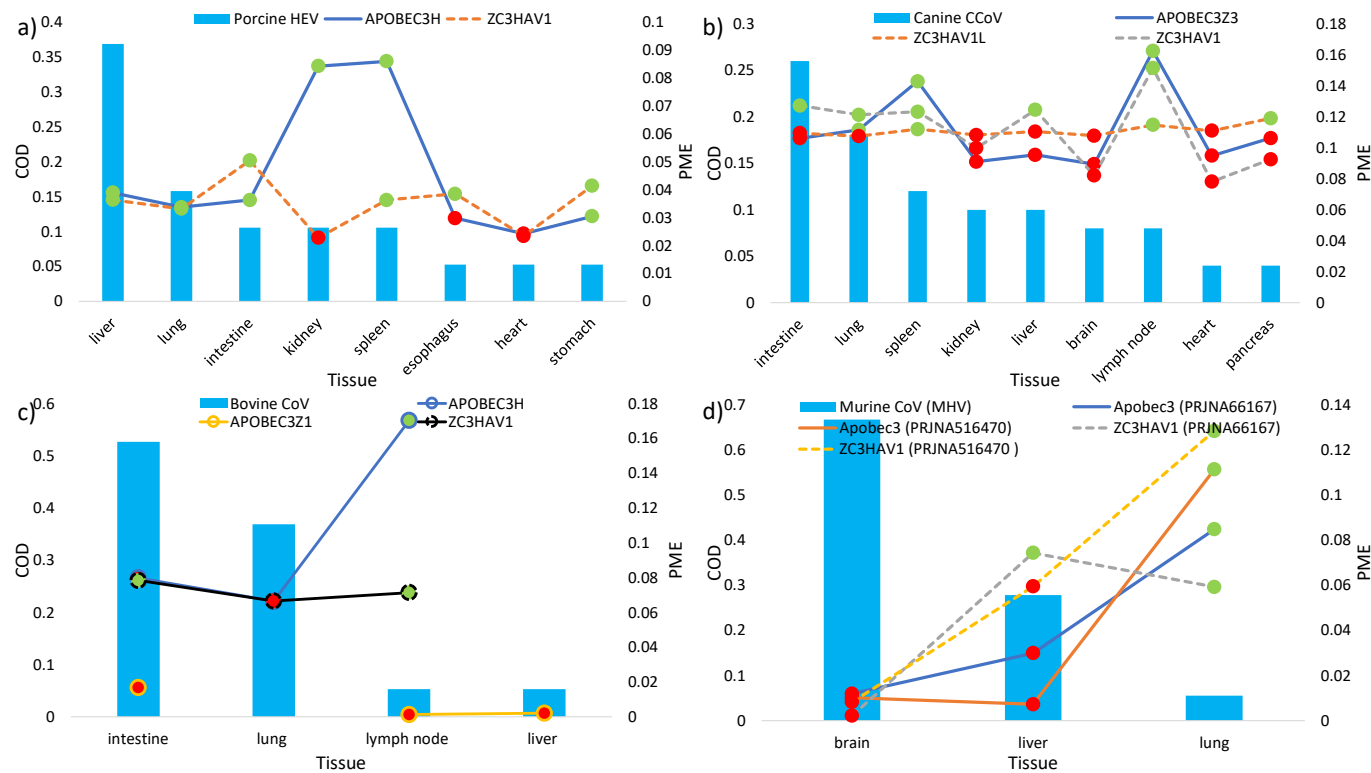
280 cattle (Fig. 2c), and Murine MHV in mice (Fig. 2d), as well as the relative mRNA expressions
281 (PMEs) for AVP isoforms in infected tissues. Like human coronaviruses, other mammalian
282 coronavirus can regularly infect tissues that exhibit both high APOBEC3 and ZAP mRNA
283 expressions, such as Porcine HEV infecting pig liver (Fig. 2a), Canine CoV infecting dog intestine
284 and lung (Fig. 2b), and Bovine CoV infecting cattle intestine (Fig. 2c). All three of these
285 coronaviruses do not avoid tissues with high AVP expressions, nor do they display a compelling
286 preference for tissues with low AVP expressions. Lastly, Murine MHV regularly infects mice
287 brain and liver but rarely infects the lung; however, mice brain and liver express low levels of
288 both APOBEC3 and ZAP PME, whereas the lung expresses high levels of both AVP PME (Fig.
289 2d). Hence, unlike the other coronaviruses, Murine MHV seems to avoid tissues with high AVP
290 expressions and prefers to infect tissues with low AVP expressions.

291 In principle, there are three possible classes of AVP expression a given tissue may conform to:
292 1) overall AVP abundance, 2) overall AVP deficiency, and 3) selective expression of AVPs. The
293 first two classes describe tissues for which both ZAP and APOBEC3 are expressed highly and
294 lowly, respectively. The third pattern can be divided into two subsets: one in which APOBEC3
295 enzymes are highly expressed and ZAP is lowly expressed, and the inverse to this pattern.
296 Figures 1 and 2 suggest that tissue-specific APOBEC3 and ZAP expressions may be correlated in
297 some species but not in others. Indeed, based on 24 human tissue, PME values of human
298 APOBEC3 and ZAP are significantly positively correlated (e.g., A3H vs ZC3HAV1: $R^2 = 0.43$, $P =$
299 0.00035). Similarly, we found significant positive correlation between the two AVPs in PME
300 values based on 17 mice tissues (ApoBec3 vs ZC3HAV1: $R^2 = 0.49$, $P = 0.0017$) and based on 10
301 dog tissues (APOBEC3Z3 vs ZC3HAV1: $R^2 = 0.56$, $P = 0.021$). However, there are no significant
302 correlation between the two AVPs in PME values based on 26 cattle tissues (APOBEC3H vs
303 ZC3HAV1: $R^2 = 0.22$, $P = 0.34$) or based on 33 pig tissues (APOBEC3H vs ZC3HAV1: $R^2 = 0.11$, $P =$
304 0.065).



305

306 Fig. 1. The histograms show the regular tissue habitats (as measured in commonness of detection COD, on primary Y axis) of SARS-
 307 CoV-2 (a, b) and of SARS-CoV and MERS (c, d). The lines represent the relative mRNA expression (in proportions of mRNA expression
 308 PME, on secondary Y axis) of a) APOBEC3 isoforms (solid lines) and b) ZAP isoforms (dash lines) in tissues susceptible to SARS-CoV-2
 309 infection, and the PME of c) APOBEC3 isoforms (solid lines) and d) ZAP isoforms (dashed lines) in tissues susceptible to SARS-CoV and
 310 MERS infections. Highlighted in green and red are PME values that are greater and lower than the averaged PME values,
 311 respectively. PME values were calculated based on averaged mRNA FPKMs retrieved from the GTEx Portal (LONSDALE *et al.* 2013).



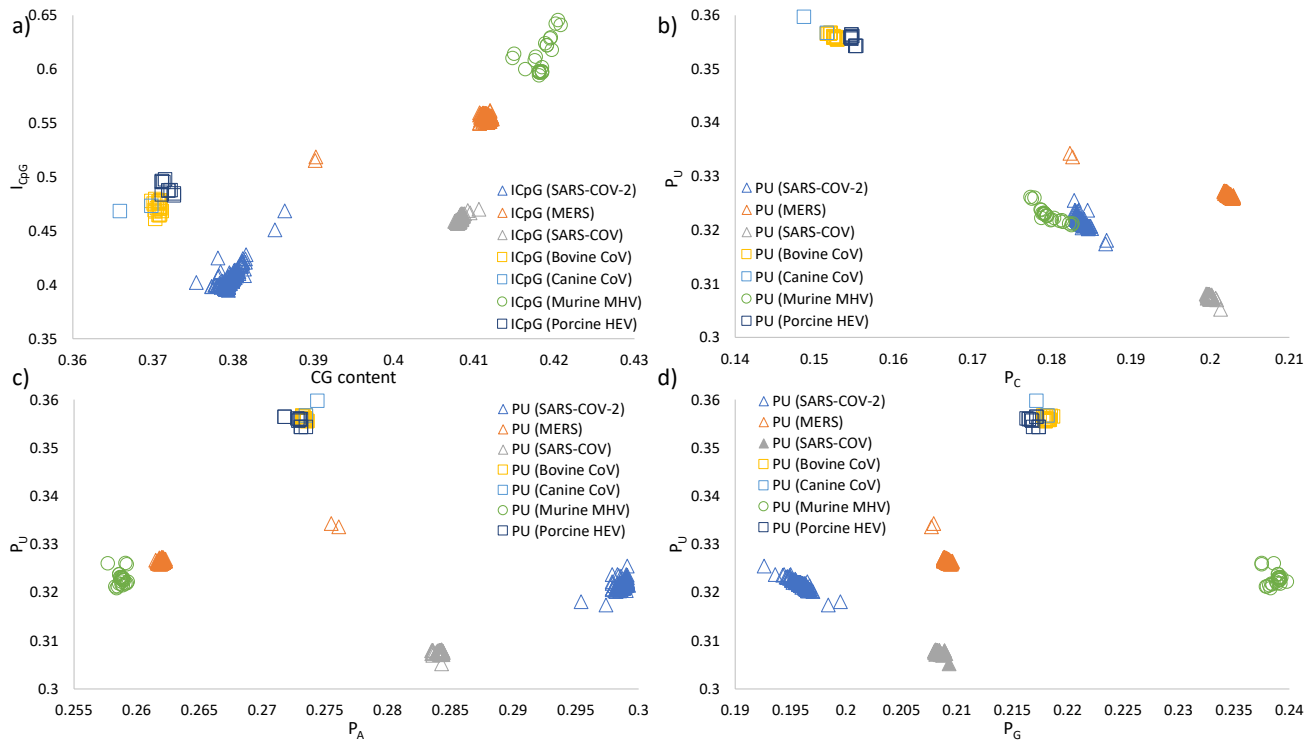
312

313 Fig. 2. The histograms show the regular tissue habitats (as measured in commonness of detection COD, on primary Y axis) of a)
 314 Porcine HEV, b) Canine CoV, c) Bovine CoV, and d) Murine MHV. The lines represent the relative mRNA expression (in proportions of
 315 mRNA expression PME, on secondary Y axis) of (a) APOBEC3 and ZAP isoforms in pig tissues susceptible to Porcine HEV infection,
 316 PME of (b) APOBEC3 and ZAP isoforms in dog tissues susceptible to Canine CoV infections, PME of (c) APOBEC3 and ZAP isoforms in
 317 cattle tissues susceptible to Bovine CoV infection, and PME of (d) APOBEC3 and ZAP isoforms in mice tissues susceptible to Murine
 318 MHV infection. Highlighted in green and red are PME values that are greater and lower than the averaged PME values, respectively.
 319 Solid lines and dashed lines represent APOBEC3 and ZAP isoforms, respectively. PME values were calculated based on averaged mRNA
 320 expressions retrieved from the Bovine Genome Database (SHAMIMUZZAMAN *et al.* 2019), BioProject PRJNA124245 (BRIGGS *et al.* 2011),
 321 TISSUE 2.0 integrated datasets (PALASCA *et al.* 2018), mouse ENCODE consortium (YUE *et al.* 2014) and BioProject PRJNA516470
 322 (NAQVI *et al.* 2019).

323 **Only coronaviruses targeting tissues with high AVP expressions exhibit decreased CpG and**
324 **increased U content**

325 In the previous section, we demonstrated that many surveyed host-specific coronaviruses
326 commonly infect tissues that exhibit high levels of AVPs (Fig. 1, 2; supplemental Fig. S1, S2), but
327 MHV does not conform to this observation (Fig. 2d). Here we compared the CpG and U content
328 of these coronaviruses and found that viruses that regularly infect AVP-rich tissues tend to
329 exhibit diminished CpG content in tandem with elevated U content. Conversely, MHV neither
330 targets AVP-rich tissues, nor does its genome indicate directional mutation with respect to CpG
331 or U content. Both trimmed genomes (Fig. 3a, see Materials and Methods) and whole genomes
332 (Supplemental Fig. S3a) shows that MHV has the highest I_{CpG} (about 0.6 or higher) while SARS-
333 CoV-2 has the lowest I_{CpG} (below 0.43 in all but two genomes). As for all other coronaviruses
334 surveyed, they also exhibit low $I_{\text{CpG}} < 0.5$ except for MERS being slightly higher. It should be
335 noted that among the 7 coronaviruses surveyed, I_{CpG} values also show the greatest variation
336 among Murine MHV genomes whereas I_{CpG} values are relatively much more constrained among
337 genomes of the other 6 (Supplemental Fig. S4a). Indeed, CpG content is weakly constrained in
338 Murine MHV.

339 Figure panels 3b, 3c, and 3d show that the proportion of U nucleotides (P_U) decreases with the
340 proportion of C nucleotides (P_C), but P_U does not correlate with P_A or P_G . This global relationship
341 in the trimmed genomes may suggest a hallmark of C to U deamination in coronaviruses, that
342 single stranded RNA genomes could indeed be subjected to editing by APOBEC3 proteins.
343 Specifically, Bovine CoV, Canine CoV, and Porcine HEV all have very high P_U and conversely very
344 low P_C . In comparison, since Murine MHV does not infect host tissues with high APOBEC3
345 expression, it may have been subjected to less C to U deamination and therefore it has notably
346 reduced P_U and increased P_C . Additionally, similar to I_{CpG} , P_U is least constrained in Murine MHV
347 in comparison to any other coronavirus (Supplemental Fig. S4b). As for human coronaviruses,
348 figure 3b shows that P_U is comparably low in SARS-CoV-2 and in MERS, like in Murine MHV, and
349 even lower in SARS-CoV. Nonetheless, it is important to note that the emergence of all three
350 human coronaviruses are much more recent in comparison to coronaviruses of other mammals;
351 their genomes had short evolutionary time to be shaped by host AVPs. Lastly, the same
352 patterns were observed when P_U was re-analyzed using whole, untrimmed, genomes
353 (Supplemental Fig. S3).



354
355

356 Fig. 3. ICpG and nucleotide proportions for seven Coronaviruses with complete genomes and
 357 host information. All genomes were aligned with MAFFT and sequence ends were trimmed (see
 358 Materials and Methods). In panel a) shows that SARS-CoV-2 has the least ICpG in comparison to
 359 other coronaviruses from their natural hosts. In panels b), c) and d) show that proportions of U
 360 (P_U) negatively correlates with P_C but not with P_A or P_G , and that P_U is similarly the highest
 361 among Bovine CoV, Canine CoV, and Porcine CoV, and similarly the lowest among Murine MHV
 362 and human coronaviruses. Each panel includes 2666 SARS-CoV-2 genomes, 403 MERS genomes,
 363 134 SARS-CoV genomes, 20 Bovine CoV genomes, two Canine CoV genomes, 26 Murine MHV
 364 genomes, and 10 Porcine HEV genomes.

365

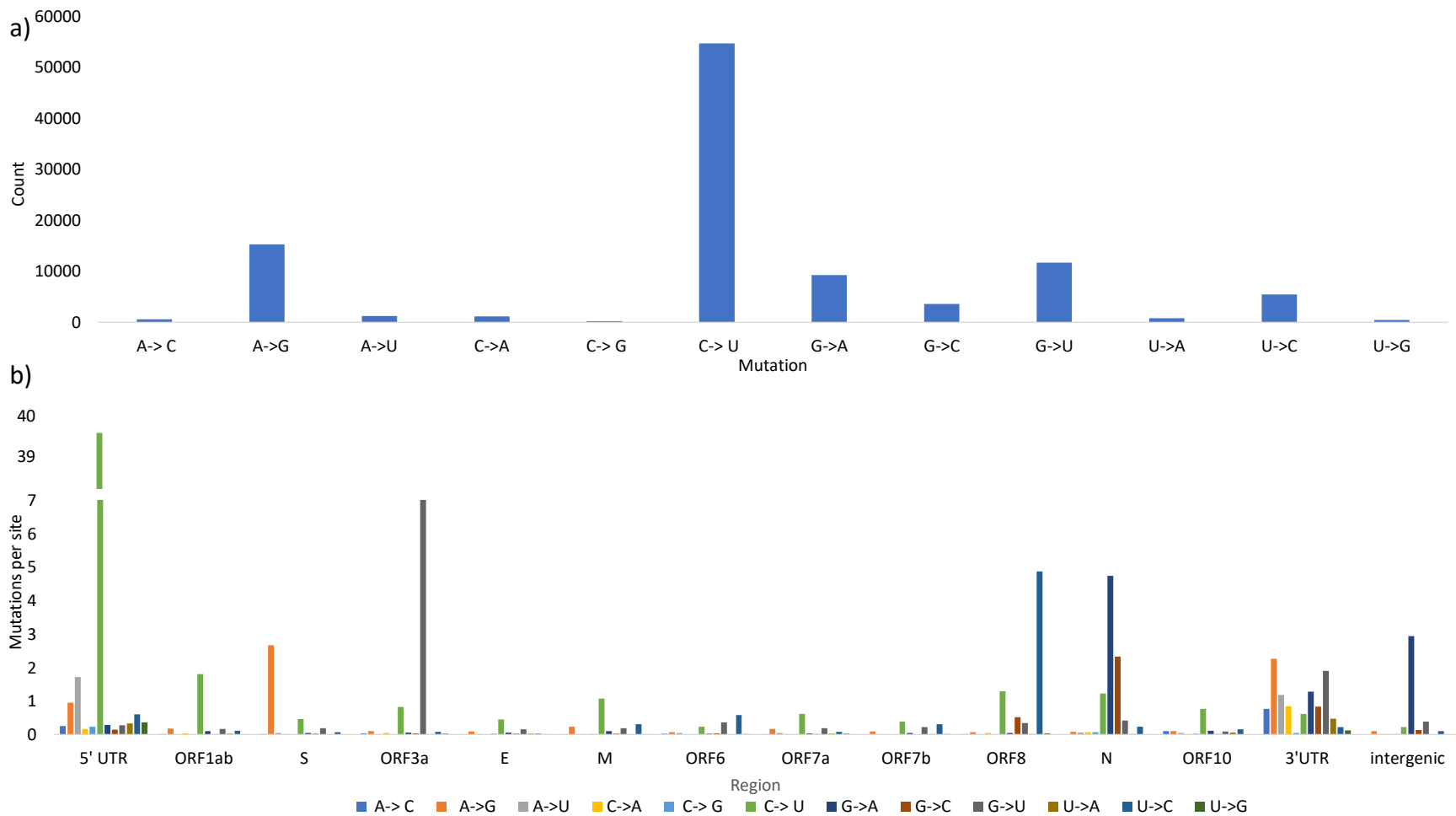
366 Strong evidence of directional mutation shaped by AVPs in local SARS-CoV-2 viral regions

367 In the above results, we demonstrated that viral genomes show more pronounced shifts
 368 towards CpG deficiency and elevated U content when the virus regularly infect host tissues with
 369 high expression of the two AVPs. However, two limitations of figure 3 are 1) it does not show
 370 the substitution patterns within local viral regions (such as the ORF1ab region), and 2) because
 371 human viruses share similar or lower global P_U in comparison to Murine MHV which
 372 predominantly infects AVP-deficient tissues, we cannot suggest that APOBEC3 would shape
 373 human coronavirus genomes over time. To resolve these limitations, we examined whether

374 there has been an evolutionary history of P_U elevation in local SARS-CoV-2 regions over the span
375 of 4 months since the virus was first isolated.

376 To resolve the first limitation, figure 4 shows single nucleotide polymorphisms (SNPs) between
377 28475 aligned SARS-CoV-2 sequences (including complete and incomplete sequences)
378 (retrieved from <https://bigd.big.ac.cn/ncov/variation/statistics>, last accessed May 16, 2020),
379 and the comparison was made against the first identified Wuhan-Hu-1 strain (MN908947).
380 Based on global sequence comparison, figure 4a shows that most SNPs are C->U substitutions.
381 More specifically, local mutation patterns (Fig. 4b) show that among 28475 sequence samples,
382 C->U substitutions are most prevalent at the 5' UTR region and the ORF1ab region (normalized
383 by region length), but not at any other viral regions. To resolve the second limitation, figure 5
384 and supplemental figure S5 show the local mutation patterns over time in a sample of 99
385 complete and high-quality SARS-CoV-2 genomes with complete NCBI annotations. Each
386 retrieved sample had been collected on a different day, since first isolation (Wuhan-Hu-1,
387 MN908947, December 31, 2019) to the most recently isolated sample (mink/NED/NB04,
388 MT457401, May 6, 2020) (see Materials and Methods), and each sample was grouped into one
389 of six time ranges. Indeed, with reference to strain Wuhan-Hu-1, aligned sequences show an
390 excessive number of C->U substitutions, and that the total number of C->U substitutions
391 increases over time, but only at the 5' UTR region (Fig. 5a) and ORF1ab region (Fig. 5b) and not
392 in other regions (Fig. 5c, 5d, Supplemental Fig. S5). It is noteworthy that in the S region,
393 directional mutation over time favours A->G substitutions; whereas in the ORF3a region,
394 directional mutation over time seems to favour G->U, but the number of G->U mutations have
395 decreased in the latest group of genome samples. Together, figures 4 and 5 suggest that
396 APOBEC3 may indeed edit single-stranded RNA genomes, specifically the 5'UTR and ORF1ab
397 regions in SARS-CoV-2. Whereas in the S region specifically, A->G directional mutation may be
398 the result of deamination by the mammalian adenosine deaminase acting on RNA type 1
399 (ADAR1) enzyme (JIANG 2020; SAMUEL 2011; ZHAO *et al.* 2004).

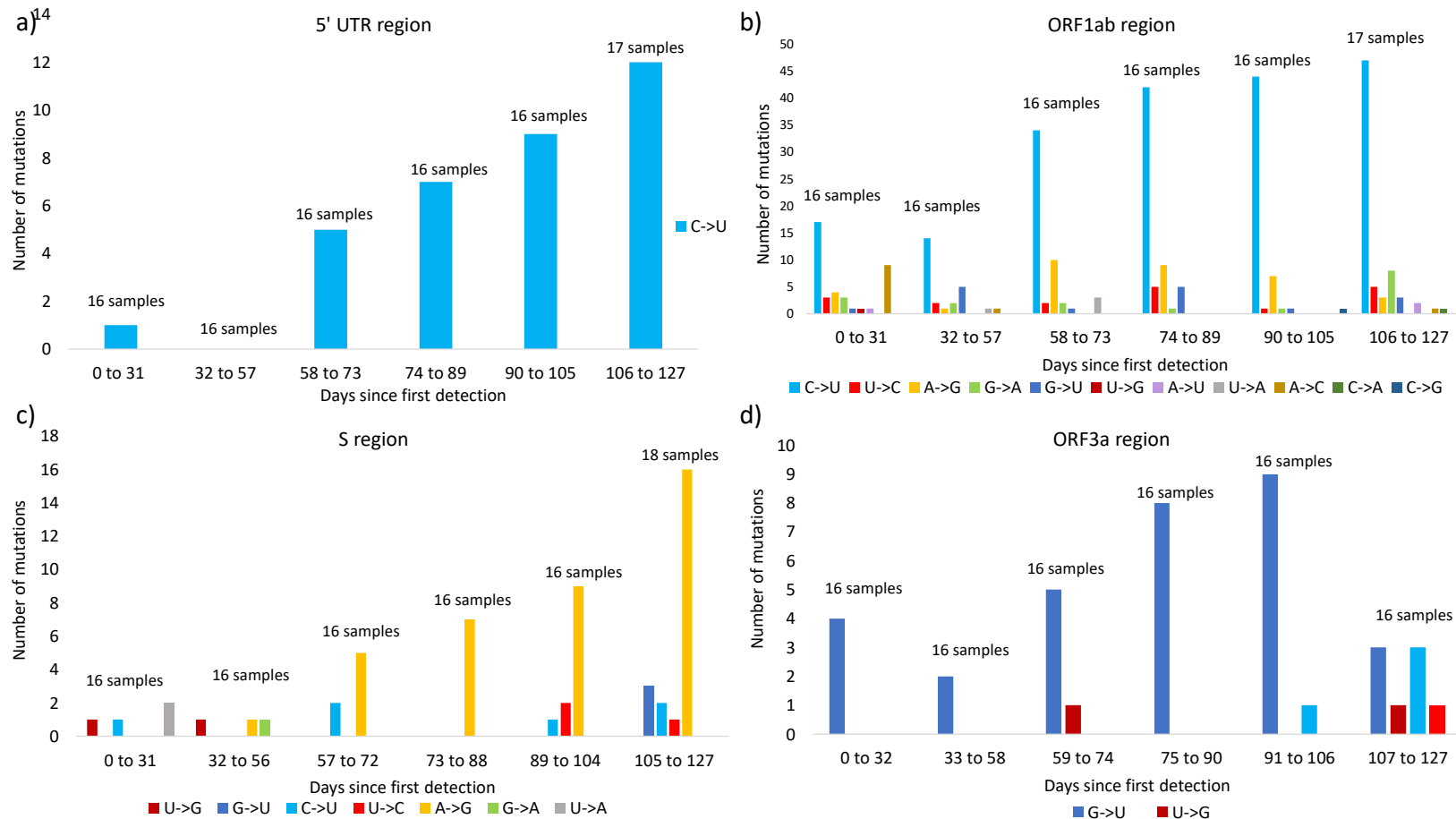
400 Lastly, we compared differences in I_{CpG} between viral regions over time among the 99 SARS-
401 CoV-2 samples. Figure 6 shows no difference in I_{CpG} between sequences sampled at different
402 time (Since December 31, 2019 to May 6, 2020). This suggests that I_{CpG} changes slowly over
403 time. However, there are notable differences in I_{CpG} among specific viral regions. In particular,
404 ORF1ab, S, and ORF6 regions have the lowest I_{CpG} values, whereas the 5' UTR, E, and ORF10
405 regions have the highest I_{CpG} values at above 1. The selective pressure for CpG deficiency to
406 evade ZAP is not uniform across different SARS-CoV-2 viral regions.



407
408

409 Fig. 4. Single nucleotide polymorphisms in 28474 SARS-CoV-2 (complete and incomplete) samples sequenced to date (May 6, 2020),
410 with reference to strain Wuhan-Hu-1 (MN908947). Panel a) shows strong global C->U directional mutation when the entire genomes
411 are considered. Panel b) shows that the number of C->U directional mutations (normalized by the length of the region) is

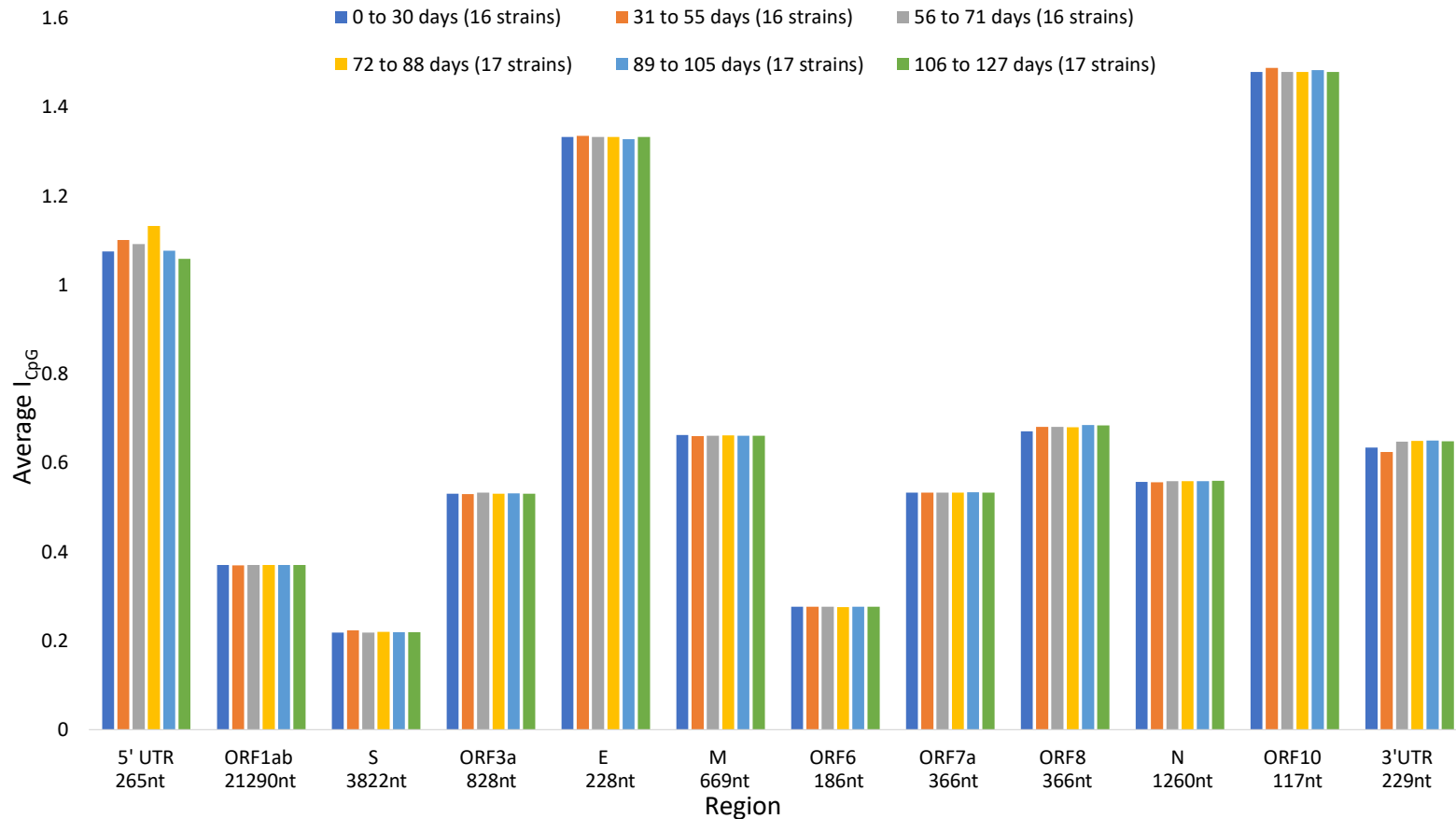
412 predominantly observed in the 5' UTR and ORF1ab viral regions but not in other regions. Indels and ambiguous point mutations
 413 were omitted.



414

415 Fig. 5. Local mutation patterns over time in a sample of 99 complete and high-quality SARS-CoV-2 sequences with complete NCBI
 416 annotations. Each retrieved sample was collected on a different day, since first isolation (Wuhan-Hu-1, MN908947, December 31,
 417 2019) to the most recent isolation (mink/NED/NB04, MT457401, May 6, 2020) (see Materials and Methods), then each sample was
 418 grouped into one of six time ranges. In panels a) and b) shows that the total number of C->U mutations are the most prevalent and

419 they increase over time, in the 5' UTR region and ORF1ab region, respectively. In panels c) and d) show that C->U mutation are not
 420 prevalent, and that A->G mutation and G->U mutations are favoured in the S and ORF3a regions, respectively.



421

422 Fig. 6. Local I_{CpG} values over time in a sample of 99 complete and high-quality SARS-CoV-2 sequences with complete NCBI
 423 annotations. Each retrieved sample was collected on a different day, since first isolation (Wuhan-Hu-1, MN908947, December 31,
 424 2019) to the most recent isolation (mink/NED/NB04, MT457401, May 6, 2020) (see Materials and Methods), then each sample was
 425 grouped into one of six time ranges. I_{CpG} does not change substantially over the 127 days since first detection, but I_{CpG} is not uniform
 426 across viral regions. I_{CpG} is the lowest in ORF1ab, S, and ORF6 regions, and the highest in the 5' UTR, E, and ORF10 regions.

427 DISCUSSION

428 SARS-CoV-2 poses a serious global health emergency. Since its outset in Wuhan City, Hubei
429 province of China in December 2019, the viral outbreak has resulted in over 7 million confirmed
430 cases around the world (<https://www.who.int/emergencies/diseases/novel-coronavirus-2019>,
431 last accessed June 11, 2020). The pandemic has prompted an immediate global effort to
432 sequence the genome of SARS-CoV-2, and over 28000 genome samples have been publicly
433 deposited over the course of just four months to facilitate vaccine development strategies.
434 With a wealth of sequence data, we performed a comprehensive comparative genome study on
435 SARS-CoV-2 and six other coronaviruses across five mammalian species, with the aim to
436 understand how coronaviruses evolve in response to tissue-specific host immune systems.

437 We tested the hypothesis that both APOBEC3 and ZAP immune responses act as primary
438 selective pressures to shape the genome of an infecting coronavirus over the course of its
439 evolutionary history within host tissues. Specifically, viral genomes are driven towards reduced
440 CpG dinucleotides to elude ZAP-mediated cellular antiviral defense, and increased U residues
441 because of RNA editing by APOBEC3 proteins. In line with our expectations, we found
442 compelling hallmarks of CpG deficiency and C to U deamination globally in mammalian
443 coronaviruses (i.e., Fig. 3: Bovine CoV, Canine CoV, and Porcine HEV) that regularly infect host
444 tissues expressing both AVPs in abundance (Fig. 2a, 2b, and 2c). Unsurprisingly, these global
445 trends were absent from Murine MHV genomes (Fig. 3) as this virus does not regularly infect
446 tissues that highly express AVPs (Fig. 2d). Corroborating this observation, both I_{CpG} and P_{U}
447 values show the greatest variation among Murine MHV genomes (Supplemental Fig. S4),
448 suggesting that MHV is not functionally constrained by either AVP. This aligns with our
449 prediction that for a virus regularly infecting host tissues that are deficient in AVPs, there will be
450 no strong directional mutations resulting in decreased CpG dinucleotides or elevated U
451 residues. Conversely, when a virus regularly infects host tissues that are abundant in ZAP and
452 APOBEC3, these AVPs shape the molecular evolution of viral genomes in two ways: CpG
453 deficiency contributes to the survival of the virus by evading ZAP-mediated antiviral defense
454 through CG dinucleotide recognition, and elevated U content as the result of genome editing by
455 APOBEC3.

456 In comparison to other mammalian coronaviruses, human coronaviruses (SARS-CoV-2, SARS-
457 COV, MERS) have been circulating in the human hosts for a much shorter time, particularly
458 SARS-CoV-2. Among the three, SARS-CoV-2 genomes shows extreme CpG deficiency (Fig. 3a); its
459 I_{CpG} values are comparable to that of the Bat CoV RaTG13 coronavirus infecting the bat species
460 *Rhinolophus affinis* (XIA 2020) but lower than that of all other coronaviruses studied herein as
461 well as all other mammalian specific coronaviruses (XIA 2020). Indeed, many recently published
462 studies point to Bat CoV RaTG13 as the most closely related known relative of SARS-CoV-2
463 when the whole genome is considered (ANDERSEN *et al.* 2020; LAI *et al.* 2020; SHANG *et al.* 2020;

464 TANG *et al.* 2020), and to *Rhinolophus affinis* as a potential intermediate host or reservoir for
465 SARS-CoV-2 (LIU *et al.* 2020). Moreover, local comparative analyses on CpG content have two
466 implications. First, SARS-CoV-2 has acquired CpG deficiency in an intermediate reservoir prior to
467 zoonotic transmission to humans, as CpG deficiency may be acquired slowly since there is no
468 notable change in I_{CpG} across all 12 viral regions in the span of four months since SARS-CoV-2
469 was initially isolated. In this context, it is regrettable that the Bat CoV RaTG13 was not
470 sequenced when it was initially sampled in 2013. The downshifting in I_{CpG} in RaTG13 would have
471 served as a warning that the virus will likely infect tissues with high ZAP expression, because the
472 viral genome has successfully evolved to evade ZAP-mediated antiviral defense in humans.
473 Second, the evolutionary pressure for CpG deficiency may be region specific. The S, ORF1ab,
474 and ORF6 regions have the most severe CpG deficiencies (Fig. 6, $I_{\text{CpG}} < 0.4$), whereas the 5' UTR,
475 E, and ORF10 regions have the highest CpG content with no signs of CpG deficiency (Fig. 6, $I_{\text{CpG}} >$
476 1). While evolution has allowed the Spike protein to elude ZAP because it is crucial for host cell
477 recognition and entry, structural genes such as the Envelope and Membrane protein are
478 subjected to less selective pressure to evade ZAP.

479 A current survey of SARS-CoV-2 genomes does not indicate drastically increased U and
480 decreased C contents. A global sequence comparison shows that SARS-CoV-2 (and SARS-CoV
481 and MERS) have comparable U and C contents as Murine MHV, but higher U and lower C
482 contents in comparison to Bovine CoV, Canine CoV, and Porcine HEV (Fig. 3b). This is because
483 while a coronavirus infecting a specific host tissue for a long time would experience the same
484 cellular antiviral environment and is consequently expected to have undergone significant RNA
485 editing; newly emerging coronaviruses such as SARS-CoV-2 would not have enough time to
486 accumulate a high number of RNA modifications. Nevertheless, global nucleotide substitution
487 patterns (Fig. 4a) show that C to U substitution is still the most prevalent among SARS-CoV-2
488 genomes collected to date. This prevalence in genome wide C to U substitutions in SARS-CoV-2
489 has been similarly reported by DI GIORGIO *et al.* (2020), who also observed the same trends in
490 SARS-CoV and to a lesser degree in MERS.

491 More importantly, local sequence comparisons among SARS-CoV-2 samples indeed show that
492 there is an evolutionary history of P_U elevation in specific SARS-CoV-2 viral regions over the
493 span of 4 months since the virus was first isolated. There is an excessive number of C to U
494 substitutions, and the prevalence of C to U mutations is increasing over time, specifically in the
495 5' UTR and ORF1ab regions (Fig. 4b, 5a, 5b). This implies that these two specific viral regions are
496 under constant C to U deamination by the APOBEC3 gene family, at least in the short term so
497 far. Another noteworthy observation is that G to A substitution is preferred in the S region and
498 the numbers of G to A substitutions are increasing over time (Fig. 5c). The preference for this
499 mutation may be caused by deamination by the mammalian adenosine deaminase acting on
500 RNA type 1 (ADAR1) enzyme (DI GIORGIO *et al.* 2020; JIANG 2020), which edits A into I, and

501 subsequently into G. Although, ADAR1 was known for targeting double-stranded RNAs, not
502 single-stranded RNA sequences (EISENBERG and LEVANON 2018; O'CONNELL *et al.* 2015; SIMMONDS
503 2020; ZHAO *et al.* 2004). Regardless, these results suggest that RNA editing by host deaminase
504 systems may indeed act on coronaviruses.

505 While it is important to determine the evolution of coronavirus genomes to understand its host
506 adaptation and specificity, this study focuses more on the evolutionary pressure and RNA
507 editing process that host immune systems exert onto viral genomes. Our aim is to prompt
508 motivations for vaccine designs in the development of attenuated pathogenic RNA viruses.
509 Previous experimental works have shown that increasing CpG dinucleotides in CpG-deficient
510 viral genomes leads to drastic decrease in viral replication and virulence (ANTZIN-ANDUETZA *et al.*
511 2017; BURNS *et al.* 2009; FROS *et al.* 2017; TRUS *et al.* 2020; TULLOCH *et al.* 2014; WASSON *et al.*
512 2017), and in recent years several studies have proposed vaccine development strategies
513 involving increased CpG to attenuate pathogenic RNA viruses (BURNS *et al.* 2009; FICARELLI *et al.*
514 2020; TRUS *et al.* 2020; TULLOCH *et al.* 2014). Among coronaviruses, SARS-CoV-2 has the most
515 extreme CpG deficiency (XIA 2020), particularly in the S protein coding region (Fig. 6). Increasing
516 CpG content at the S protein may provide a good starting point for strategies to inhibit SARS-
517 CoV-2's ability to recognize and enter host cells. On the other hand, because C to U
518 deamination cannot be proof-read by viral exonuclease Nsp14-ExoN (ECKERLE *et al.* 2010; SMITH
519 *et al.* 2013; VICTOROVICH *et al.* 2020), host innate deaminases may drive up the rate of evolution
520 in viral genomes (DI GIORGIO *et al.* 2020) or modify CpG into UpG to further increase CpG
521 deficiency and reduce viral susceptibility by ZAP. The possibility of APOBEC3 editing activity
522 acting on RNA viruses and its potential exploits by viruses such as SARS-CoV-2 in the long term
523 require further investigation and scrutiny.

524 **ACKNOWLEDGEMENTS**

525 This work is supported by the Natural Sciences and Engineering Research Council of Canada
526 (NSERC) Discovery Grant to X.X. [RGPIN/2018-03878], and NSERC Doctoral Scholarship to Y.W.
527 [CGSD/2019-535291].

528 **AUTHOR CONTRIBUTIONS**

529 Y.W. and X.X. designed the study. Y.W., J.R.S., and X.X. wrote the manuscript. Y.W., P.A. and
530 X.X. collected the data. Y.W. and J. R. S. analyzed the data. Y.W., P.A., and J. R. S. prepared all
531 figures. All authors reviewed the manuscript. X.X. supervised the project.

532 **COMPETING INTERESTS**

533 The authors declare no competing interests.

534 **REFERENCES**

- 535 ANDERSEN, K. G., A. RAMBAUT, W. I. LIPKIN, E. C. HOLMES and R. F. GARRY, 2020 The proximal origin of
536 SARS-CoV-2. *Nature Medicine* **26**: 450-452.
- 537 ANTZIN-ANDUETZA, I., C. MAHIET, L. A. GRANGER, C. ODENDALL and C. M. SWANSON, 2017 Increasing the
538 CpG dinucleotide abundance in the HIV-1 genomic RNA inhibits viral replication.
539 *Retrovirology* **14**: 017-0374.
- 540 ATKINSON, N. J., J. WITTEVELDT, D. J. EVANS and P. SIMMONDS, 2014 The influence of CpG and UpA
541 dinucleotide frequencies on RNA virus replication and characterization of the innate
542 cellular pathways underlying virus attenuation and enhanced replication. *Nucleic Acids*
543 *Res* **42**: 4527-4545.
- 544 BISHOP, K. N., R. K. HOLMES, A. M. SHEEHY and M. H. MALIM, 2004 APOBEC-mediated editing of viral
545 RNA. *Science* **305**: 1100658.
- 546 BOGERD, H. P., and B. R. CULLEN, 2008 Single-stranded RNA facilitates nucleocapsid: APOBEC3G
547 complex formation. *RNA (New York, N.Y.)* **14**: 1228-1236.
- 548 BRIGGS, J., M. PAOLONI, Q. R. CHEN, X. WEN, J. KHAN *et al.*, 2011 A compendium of canine normal
549 tissue gene expression. *PLoS One* **6**: 31.
- 550 BURNS, C. C., R. CAMPAGNOLI, J. SHAW, A. VINCENT, J. JORBA *et al.*, 2009 Genetic Inactivation of
551 Poliovirus Infectivity by Increasing the Frequencies of CpG and UpA Dinucleotides within
552 and across Synonymous Capsid Region Codons. *Journal of Virology* **83**: 9957.
- 553 CARDON, L. R., C. BURGE, D. A. CLAYTON and S. KARLIN, 1994 Pervasive CpG suppression in animal
554 mitochondrial genomes. *Proc Natl Acad Sci U S A* **91**: 3799-3803.
- 555 CHIU, Y.-L., and W. C. GREENE, 2008 The APOBEC3 Cytidine Deaminases: An Innate Defensive
556 Network Opposing Exogenous Retroviruses and Endogenous Retroelements. *Annual*
557 *Review of Immunology* **26**: 317-353.
- 558 CULLEN, B. R., 2006 Role and Mechanism of Action of the APOBEC3 Family of Antiretroviral
559 Resistance Factors. *Journal of Virology* **80**: 1067.
- 560 DI GIOACCHINO, A., P. ŠULC, A. V. KOMAROVA, B. D. GREENBAUM, R. MONASSON *et al.*, 2020 The
561 heterogeneous landscape and early evolution of pathogen-associated CpG and UpA
562 dinucleotides in SARS-CoV-2. *bioRxiv*: 2020.2005.2006.074039.
- 563 DI GIORGIO, S., F. MARTIGNANO, M. G. TORCIA, G. MATTIUZ and S. G. CONTICELLO, 2020 Evidence for
564 host-dependent RNA editing in the transcriptome of SARS-CoV-2. *Science Advances*:
565 eabb5813.
- 566 DUFF, M. O., S. OLSON, X. WEI, S. C. GARRETT, A. OSMAN *et al.*, 2015 Genome-wide identification of
567 zero nucleotide recursive splicing in *Drosophila*. *Nature* **521**: 376-379.
- 568 ECKERLE, L. D., M. M. BECKER, R. A. HALPIN, K. LI, E. VENTER *et al.*, 2010 Infidelity of SARS-CoV
569 Nsp14-exonuclease mutant virus replication is revealed by complete genome
570 sequencing. *PLoS Pathog* **6**: 1000896.

- 571 EISENBERG, E., and E. Y. LEVANON, 2018 A-to-I RNA editing — immune protector and transcriptome
572 diversifier. *Nature Reviews Genetics* **19**: 473-490.
- 573 FAGERBERG, L., B. M. HALLSTRÖM, P. OKSVOLD, C. KAMPF, D. DJUREINOVIC *et al.*, 2014 Analysis of the
574 Human Tissue-specific Expression by Genome-wide Integration of Transcriptomics and
575 Antibody-based Proteomics. *Molecular & Cellular Proteomics* **13**: 397.
- 576 FICARELLI, M., I. ANTZIN-ANDUETZA, R. HUGH-WHITE, A. E. FIRTH, H. SERTKAYA *et al.*, 2020 CpG
577 Dinucleotides Inhibit HIV-1 Replication through Zinc Finger Antiviral Protein (ZAP)-
578 Dependent and -Independent Mechanisms. *J Virol* **94**: 01337-01319.
- 579 FROS, J. J., I. DIETRICH, K. ALSHAIKHAHMED, T. C. PASSCHIER, D. J. EVANS *et al.*, 2017 CpG and UpA
580 dinucleotides in both coding and non-coding regions of echovirus 7 inhibit replication
581 initiation post-entry. *Elife* **29**: 29112.
- 582 GREENBAUM, B. D., A. J. LEVINE, G. BHANOT and R. RABADAN, 2008 Patterns of evolution and host
583 gene mimicry in influenza and other RNA viruses. *PLoS Pathog* **4**: 1000079.
- 584 GREENBAUM, B. D., R. RABADAN and A. J. LEVINE, 2009 Patterns of oligonucleotide sequences in viral
585 and host cell RNA identify mediators of the host innate immune system. *PLoS One* **4**:
586 0005969.
- 587 GUO, X., J. MA, J. SUN and G. GAO, 2007 The zinc-finger antiviral protein recruits the RNA
588 processing exosome to degrade the target mRNA. *Proceedings of the National Academy
589 of Sciences* **104**: 151.
- 590 HARRIS, R. S., K. N. BISHOP, A. M. SHEEHY, H. M. CRAIG, S. K. PETERSEN-MAHRT *et al.*, 2003 DNA
591 deamination mediates innate immunity to retroviral infection. *Cell* **113**: 803-809.
- 592 HARRIS, R. S., and J. P. DUDLEY, 2015 APOBECs and virus restriction. *Virology* **480**: 131-145.
- 593 HAYWARD, J. A., M. TACHEDJIAN, J. CUI, A. Z. CHENG, A. JOHNSON *et al.*, 2018 Differential Evolution of
594 Antiretroviral Restriction Factors in Pteropid Bats as Revealed by APOBEC3 Gene
595 Complexity. *Molecular Biology and Evolution* **35**: 1626-1637.
- 596 JIANG, W., 2020 Mutation Profile of Over 4,500 SARS-CoV-2 Isolations Reveals Prevalent
597 Cytosine-to-Uridine Deamination on Viral RNAs. Preprints.
- 598 KARLIN, S., J. MRÁZEK and A. M. CAMPBELL, 1997 Compositional biases of bacterial genomes and
599 evolutionary implications. *J Bacteriol* **179**: 3899-3913.
- 600 KATO, K., and D. M. STANDLEY, 2013 MAFFT multiple sequence alignment software version 7:
601 improvements in performance and usability. *Molecular Biology and Evolution* **30**: 772-
602 780.
- 603 KIM, D., J.-Y. LEE, J.-S. YANG, J. W. KIM, V. N. KIM *et al.*, 2020 The Architecture of SARS-CoV-2
604 Transcriptome. *Cell* **181**: 914-921.e910.
- 605 LAI, C.-C., T.-P. SHIH, W.-C. KO, H.-J. TANG and P.-R. HSUEH, 2020 Severe acute respiratory
606 syndrome coronavirus 2 (SARS-CoV-2) and coronavirus disease-2019 (COVID-19): The
607 epidemic and the challenges. *International Journal of Antimicrobial Agents* **55**: 105924.

- 608 LIU, P., J. Z. JIANG, X. F. WAN, Y. HUA, L. LI *et al.*, 2020 Are pangolins the intermediate host of the
609 2019 novel coronavirus (SARS-CoV-2)? PLoS Pathog **16**.
- 610 LONSDALE, J., J. THOMAS, M. SALVATORE, R. PHILLIPS, E. LO *et al.*, 2013 The Genotype-Tissue
611 Expression (GTEx) project. Nature Genetics **45**: 580-585.
- 612 MANGEAT, B., P. TURELLI, G. CARON, M. FRIEDLI, L. PERRIN *et al.*, 2003 Broad antiretroviral defence by
613 human APOBEC3G through lethal editing of nascent reverse transcripts. Nature **424**: 99-
614 103.
- 615 MEAGHER, J. L., M. TAKATA, D. GONÇALVES-CARNEIRO, S. C. KEANE, A. REBENDENNE *et al.*, 2019 Structure
616 of the zinc-finger antiviral protein in complex with RNA reveals a mechanism for
617 selective targeting of CG-rich viral sequences. Proc Natl Acad Sci U S A **116**: 24303-
618 24309.
- 619 MILEWSKA, A., E. KINDLER, P. VKOVSKI, S. ZEGLEN, M. OCHMAN *et al.*, 2018 APOBEC3-mediated
620 restriction of RNA virus replication. Scientific Reports **8**: 5960.
- 621 NABEL, C. S., J. W. LEE, L. C. WANG and R. M. KOHLI, 2013 Nucleic acid determinants for selective
622 deamination of DNA over RNA by activation-induced deaminase. Proceedings of the
623 National Academy of Sciences: 201306345.
- 624 NAQVI, S., A. K. GODFREY, J. F. HUGHES, M. L. GOODHEART, R. N. MITCHELL *et al.*, 2019 Conservation,
625 acquisition, and functional impact of sex-biased gene expression in mammals. Science
626 **365**.
- 627 O'CONNELL, M. A., N. M. MANNION and L. P. KEEGAN, 2015 The Epitranscriptome and Innate
628 Immunity. PLoS genetics **11**: e1005687-e1005687.
- 629 ODON, V., J. J. FROS, N. GOONAWARDANE, I. DIETRICH, A. IBRAHIM *et al.*, 2019 The role of ZAP and
630 OAS3/RNaseL pathways in the attenuation of an RNA virus with elevated frequencies of
631 CpG and UpA dinucleotides. Nucleic Acids Research **47**: 8061-8083.
- 632 PALASCA, O., A. SANTOS, C. STOLTE, J. GORODKIN and L. J. JENSEN, 2018 TISSUES 2.0: an integrative
633 web resource on mammalian tissue expression. Database **2018**.
- 634 RODRIGUEZ-FRIAS, F., M. BUTI, D. TABERNEIRO and M. HOMS, 2013 Quasispecies structure,
635 cornerstone of hepatitis B virus infection: mass sequencing approach. World journal of
636 gastroenterology **19**: 6995-7023.
- 637 SAMUEL, C., 2011 Adenosine Deaminases Acting on RNA (ADARs) are Both Antiviral and Proviral.
638 Virology **411**: 180-193.
- 639 SHAMIMUZZAMAN, M., J. J. LE TOURNEAU, D. R. UNNI, C. M. DIESH, D. A. TRIANT *et al.*, 2019 Bovine
640 Genome Database: new annotation tools for a new reference genome. Nucleic Acids
641 Research **48**: D676-D681.
- 642 SHANG, J., G. YE, K. SHI, Y. WAN, C. LUO *et al.*, 2020 Structural basis of receptor recognition by
643 SARS-CoV-2. Nature **581**: 221-224.
- 644 SHARMA, S., S. K. PATNAIK, R. T. TAGGART and B. E. BAYSAL, 2016 The double-domain cytidine
645 deaminase APOBEC3G is a cellular site-specific RNA editing enzyme. Sci Rep **6**.

- 646 SHARMA, S., S. K. PATNAIK, R. THOMAS TAGGART, E. D. KANNISTO, S. M. ENRIQUEZ *et al.*, 2015 APOBEC3A
647 cytidine deaminase induces RNA editing in monocytes and macrophages. *Nature*
648 *Communications* **6**: 6881.
- 649 SHARMA, S., J. WANG, E. ALQASSIM, S. PORTWOOD, E. CORTES GOMEZ *et al.*, 2019 Mitochondrial
650 hypoxic stress induces widespread RNA editing by APOBEC3G in natural killer cells.
651 *Genome Biol* **20**: 019-1651.
- 652 SHEEHY, A. M., N. C. GADDIS, J. D. CHOI and M. H. MALIM, 2002 Isolation of a human gene that
653 inhibits HIV-1 infection and is suppressed by the viral Vif protein. *Nature* **418**: 646-650.
- 654 SIMMONDS, P., 2020 Rampant C->U hypermutation in the genomes of SARS-CoV-2 and
655 other coronaviruses – causes and consequences for their short and long evolutionary
656 trajectories. *bioRxiv*: 2020.2005.2001.072330.
- 657 SMITH, E. C., H. BLANC, M. C. SURDEL, M. VIGNUZZI and M. R. DENISON, 2013 Coronaviruses lacking
658 exoribonuclease activity are susceptible to lethal mutagenesis: evidence for
659 proofreading and potential therapeutics. *PLoS Pathog* **9**: 15.
- 660 TAKATA, M. A., D. GONÇALVES-CARNEIRO, T. M. ZANG, S. J. SOLL, A. YORK *et al.*, 2017 CG dinucleotide
661 suppression enables antiviral defence targeting non-self RNA. *Nature* **550**: 124-127.
- 662 TANG, X., C. WU, X. LI, Y. SONG, X. YAO *et al.*, 2020 On the origin and continuing evolution of SARS-
663 CoV-2. *National Science Review*.
- 664 THEYS, K., A. F. FEDER, M. GELBART, M. HARTL, A. STERN *et al.*, 2018 Within-patient mutation
665 frequencies reveal fitness costs of CpG dinucleotides and drastic amino acid changes in
666 HIV. *PLoS genetics* **14**: e1007420-e1007420.
- 667 TRUS, I., D. UDENZE, N. BERUBE, C. WHELER, M.-J. MARTEL *et al.*, 2020 CpG-Recoding in Zika Virus
668 Genome Causes Host-Age-Dependent Attenuation of Infection With Protection Against
669 Lethal Heterologous Challenge in Mice. *Frontiers in immunology* **10**: 3077-3077.
- 670 TULLOCH, F., N. J. ATKINSON, D. J. EVANS, M. D. RYAN and P. SIMMONDS, 2014 RNA virus attenuation
671 by codon pair deoptimisation is an artefact of increases in CpG/UpA dinucleotide
672 frequencies. *Elife* **3**: e04531-e04531.
- 673 VICTOROVICH, K. V., G. RAJANISH, K. T. ALEKSANDROVNA, K. S. KRISHNA, S. A. NICOLAEVICH *et al.*, 2020
674 Translation-associated mutational U-pressure in the first ORF of SARS-CoV-2 and other
675 coronaviruses. *bioRxiv*: 2020.2005.2005.078238.
- 676 WANG, S. M., and C. T. WANG, 2009 APOBEC3G cytidine deaminase association with coronavirus
677 nucleocapsid protein. *Virology* **388**: 112-120.
- 678 WASSON, M. K., J. BORKAKOTI, A. KUMAR, B. BISWAS and P. VIVEKANANDAN, 2017 The CpG dinucleotide
679 content of the HIV-1 envelope gene may predict disease progression. *Scientific Reports*
680 **7**: 8162-8162.
- 681 XIA, X., 2018 DAMBE7: New and Improved Tools for Data Analysis in Molecular Biology and
682 Evolution. *Molecular Biology and Evolution* **35**: 1550-1552.

- 683 XIA, X., 2020 Extreme Genomic CpG Deficiency in SARS-CoV-2 and Evasion of Host Antiviral
684 Defense. *Molecular Biology and Evolution*.
- 685 YAP, Y. L., X. W. ZHANG and A. DANCHIN, 2003 Relationship of SARS-CoV to other pathogenic RNA
686 viruses explored by tetranucleotide usage profiling. *BMC Bioinformatics* **4**: 1471-2105.
- 687 YUE, F., Y. CHENG, A. BRESCHI, J. VIERSTRA, W. WU *et al.*, 2014 A comparative encyclopedia of DNA
688 elements in the mouse genome. *Nature* **515**: 355-364.
- 689 ZHANG, W., J. DU, K. YU, T. WANG, X. YONG *et al.*, 2010 Association of Potent Human Antiviral
690 Cytidine Deaminases with 7SL RNA and Viral RNP in HIV-1 Virions. *Journal of Virology* **84**:
691 12903.
- 692 ZHAO, Z., H. LI, X. WU, Y. ZHONG, K. ZHANG *et al.*, 2004 Moderate mutation rate in the SARS
693 coronavirus genome and its implications. *BMC Evolutionary Biology* **4**: 21.
- 694 ZHENG, Y.-H., D. IRWIN, T. KUROSU, K. TOKUNAGA, T. SATA *et al.*, 2004 Human APOBEC3F Is Another
695 Host Factor That Blocks Human Immunodeficiency Virus Type 1 Replication. *Journal of*
696 *Virology* **78**: 6073.
- 697 ZHU, Y., G. CHEN, F. LV, X. WANG, X. JI *et al.*, 2011 Zinc-finger antiviral protein inhibits HIV-1
698 infection by selectively targeting multiply spliced viral mRNAs for degradation.
699 *Proceedings of the National Academy of Sciences* **108**: 15834.

700

701