

1 **Putative host-derived insertions in the genomes of**
2 **circulating SARS-CoV-2 variants**

3 Yiyang Yang¹, Keith Dufault-Thompson¹, Rafaela Salgado Fontenele¹, Xiaofang
4 Jiang^{1,*}

5

6 ¹National Library of Medicine, National Institutes of Health, Bethesda, Maryland,
7 USA

8 * Correspondence: Xiaofang Jiang (xiaofang.jiang@nih.gov)

9

10 **ABSTRACT**

11 Insertions in the SARS-CoV-2 genome have the potential to drive viral evolution, but
12 the source of the insertions is often unknown. Recent proposals have suggested that
13 human RNAs could be a source of some insertions, but the small size of many
14 insertions makes this difficult to confirm. Through an analysis of available direct
15 RNA sequencing data from SARS-CoV-2 infected cells, we show that viral-host
16 chimeric RNAs are formed through what are likely stochastic RNA-dependent RNA
17 polymerase template switching events. Through an analysis of the publicly available
18 GISAID SARS-CoV-2 genome collection, we identified two genomic insertions in
19 circulating SARS-CoV-2 variants that are identical to regions of the human 18S and
20 28S rRNAs. These results provide direct evidence of the formation of viral-host
21 chimeric sequences and the integration of host genetic material into the SARS-CoV-2
22 genome, highlighting the potential importance of host-derived insertions in viral
23 evolution.

24 **IMPORTANCE**

25 Throughout the COVID-19 pandemic, the sequencing of SARS-CoV-2 genomes has
26 revealed the presence of insertions in multiple globally circulating lineages of
27 SARS-CoV-2, including the Omicron variant. The human genome has been suggested
28 to be the source of some of the larger insertions, but evidence for this kind of event
29 occurring is still lacking. Here, we leverage direct RNA sequencing data and

30 SARS-CoV-2 genomes to show host-viral chimeric RNAs are generated in infected
31 cells and two large genomic insertions have likely been formed through the
32 incorporation of host rRNA fragments into the SARS-CoV-2 genome. These
33 host-derived insertions may increase the genetic diversity of SARS-CoV-2 and
34 expand its strategies to acquire genetic materials, potentially enhancing its
35 adaptability, virulence, and spread.

36 **KEYWORDS**

37 SARS-CoV-2, insertion, host-virus chimeric reads, ribosomal RNA

38

39 **INTRODUCTION**

40 During the COVID-19 pandemic, insertions have been frequently acquired in
41 SARS-CoV-2 lineages (1-4). Insertions have been associated with several globally
42 circulating lineages, including the insertion of one amino acid at position 146 of the S
43 protein (ins146N) of the variant of interest Mu (B.1.621) (4), insertions at the
44 recurrent insertion site 214 of the NTD region on the S protein that occurred in the
45 lineages B.1.214.2 (ins214TDR) and A.2.5 (ins214AAG) (1), and the insertion
46 ins214EPE in the recently-emerged variant of concern Omicron (5). Although there is
47 insufficient evidence to show the direct impact these insertions have on viral spread
48 and interference with immune responses, the fact that variants carrying those
49 insertions have circulated for long periods suggests that they might be advantageous

50 or neutral for the transmission. Results from a long-term *in vitro* experiment where
51 SARS-CoV-2 was co-incubated with highly neutralizing antibodies have also shown
52 that an 11 amino acid insertion (ins248KTRNKSTSRRE) at the NTD N5 loop of the
53 S protein was able to drive antibody escape suggesting a potential role of insertions in
54 enhancing infectivity and virulence (6). Taken together, insertions have the potential
55 to increase genetic diversity in SARS-CoV-2 and contribute to the continued
56 evolution of the virus.

57 Previous research has shown that most small insertions in the SARS-CoV-2 genome
58 likely originated from template sliding, local duplication, or template switching
59 between viruses (2). Longer insertions (equal or larger than nine nucleotides) have
60 been detected in multiple coronavirus genomes, including in variants of concern like
61 the Omicron variant, but their origin remains unknown. Host genetic material has
62 been suggested as a possible source for these insertions (5, 7). Venkatakrishnan et al.
63 suggested that the unique insertion (ins214EPE) in the Omicron variant could have
64 originated from the human common cold virus HCoV-229E or the human genome
65 based on BLAST search (5), and the human genome has been speculated to be the
66 source of multiple other small insertions (7). However, given that these insertion
67 sequences are typically short, sequence comparisons tend to be less informative, and
68 false-positive matches have a high chance of occurring. Additionally, coronavirus
69 replication occurs in modified endoplasmic reticulum-derived double-membrane
70 vesicles, providing a physical barrier between viral and host genetic material (8), and

71 coronavirus replication complexes are known to contain enzymes with proofreading
72 activity (9), both of which likely play roles in limiting the formation of host-virus
73 chimeric sequences.

74 Human-derived insertions in the SARS-CoV-2 genome would likely be generated
75 through RdRp-driven template switching events between SARS-CoV-2 and host
76 mRNA. While template switching events between coronaviruses are common (10-13)
77 and likely contribute to the emergence of SARS-CoV-2 lineages including the
78 deltacron variant (14), template switching events between coronaviruses and host
79 RNAs are rarely documented (15, 16). Chimeric reads between SARS-CoV-2 RNA
80 and human RNA have been detected but were interpreted as a signal of SARS-CoV-2
81 integration into the human genome in a previous controversial study (17). Others have
82 suggested that the chimeric reads were likely to be template switching artifacts
83 mediated by reverse transcriptase or PCR during library preparation (18-21). One
84 possible explanation that was largely omitted in these studies is that the
85 SARS-CoV-2-host chimeric RNA could be generated by RdRp-driven template
86 switching.

87 Here, to investigate the possible existence of SARS-CoV-2-host chimeric RNA, we
88 take advantage of the publically available Nanopore direct RNA sequencing data of
89 SARS-CoV-2. Direct RNA-seq sequences the individual polyadenylated RNAs
90 directly mitigating the possible formation of chimeric reads during library preparation
91 or amplification. We first identified SARS-CoV-2-host chimeric RNA from direct

92 RNA-seq data and showed that RdRp-driven template switching between
93 SARS-CoV-2 and host mRNA occurs, but it is infrequent and stochastic. We also
94 found that highly expressed host genes and structural RNA genes have a higher
95 chance to be observed in chimeric RNA reads. We then systematically analyzed the
96 SARS-CoV-2 genomes deposited in the GISAID database (22), resulting in the
97 identification of two insertions in functional SARS-CoV-2 genomes that likely
98 originated from the host 18S and 28S rRNAs.

99 **RESULTS**

100 **Host-virus mRNA chimera are rare but do exist**

101 We first analyzed direct RNA-seq data from SARS-CoV-2 infected cell lines to
102 identify sequences formed from chimeric host-viral RNAs. The direct RNA-seq data
103 were quality filtered and mapped to both the host and SARS-CoV-2 transcriptomes to
104 identify potential chimeric sequences. Out of the 30 samples that were analyzed,
105 host-viral chimeric reads were detected in 16 of the samples with an average of 0.029%
106 (standard deviation 0.048%) of the reads mapped to SARS-CoV-2 being chimeric
107 (Supplementary Table 1). Chimeric reads were typically rare, making up 0.206% of
108 one sample, but less than 0.06% of the other 15 samples, and these rates may be an
109 overestimation due to the cell lines used compared to what would be observed in *in*
110 *vivo* conditions. Additionally, chimeric reads detected in five samples were further
111 investigated using paired-end sequencing short reads from the same samples
112 (Supplementary Table 2). Approximately 1.4% (5 out of 357) of chimeric reads were

113 supported by at least five read pairs spanning the junctions. This finding implies that a
114 small fraction of the host-viral chimeric mRNA molecules could function as templates
115 for RNA replication.

116 We then analyzed the chimeric reads to identify trends in how the viral and host RNA
117 sequences were joined. All the viral-derived sequences in chimeric reads were
118 annotated as positive-sense RNA and a majority (92.49%) of the reads contained
119 host-derived positive-sense sequences. Upon further examination, the few host reads
120 that were identified as being negative-sense were largely long non-coding RNAs that
121 were present in the raw reads as the negative-sense sequences, making it likely that
122 they were mis-annotated rather than actually being derived from negative-sense RNA.
123 These results suggest that the host-viral chimeric sequences are not the result of the
124 integration of the viral genetic material into the host genome, which would have
125 resulted in a nearly equal mix of positive and negative sense viral sequences (17).
126 Most likely, these host-viral chimeric sequences were created from
127 positive-to-positive-strand template switching events (23, 24).

128 **Viral-host chimeric read formation is likely a stochastic process**

129 The chimeric reads were then analyzed to determine if there were any patterns in the
130 composition of the sequences and in which positions relative to the references they
131 were formed. Both viral to host and host to viral chimeric sequences were detected in
132 the direct RNA-seq data, but the chimeric reads did not show a preference for either
133 organization (Supplementary Table 1). Both types of sequences were seen in

134 approximately the same frequency, with viral to host reads making up 55% of the
135 chimeric sequences and host to viral reads making up 45%. This lack of strong
136 preference may indicate that host RNA can be readily recognized by viral RdRp, but
137 other factors like the exclusion of host RNA by the formation of the
138 double-membrane vesicles might prevent the formation of chimeric RNAs. When
139 examining the positions of the junctions on the viral RNA sequences, we found there
140 was a bias toward the junction sites being located in the dense coding region near the
141 three prime end of the sequence, with fewer junctions being identified in the ORF1ab
142 genes, the largest region of the genome (Fig.1). This is likely due to the ORF1ab
143 region not being retained in the canonical SARS-CoV-2 subgenomic RNAs resulting
144 in fewer viral RNAs being synthesized with these regions that could form chimeric
145 RNAs (25). It suggests that the process by which chimeric sequences are formed is
146 likely stochastic, depending on the availability of template RNA molecules.

147 Previous studies have also found that indel formation and template switching events
148 preferentially occur in the loops and stems formed in the RNA secondary structure (2,
149 3). First a permutation test was used to investigate if junction sites were commonly
150 located in stems (positions that form base-pairs) or non-stem regions (non-base-paired
151 positions) in the viral RNA. The results of this test showed a significant (P-values <
152 0.01) preference for the formation of junctions in non-base-paired regions of the RNA
153 secondary structure (Fig.1). One-sided Fisher's exact tests were performed to explore
154 if junction sites were enriched in specific types of RNA structures. Consistent with the

155 results of the permutation test, stems were under-represented at the junction sites
156 (Supplementary Table 3). We speculate that the non-base-paired regions of the
157 SARS-CoV-2 RNA may be more susceptible to stochastic template-switching events
158 due to their more “open” configurations, where the viral RdRp could easily attach or
159 detach as it moves along the RNA.

160 An examination of the types of human gene sequences found in the chimeric
161 sequences revealed an enrichment of non-coding RNAs and highly expressed genes.
162 We found that a disproportionate number of non-coding RNAs, mainly long
163 non-coding RNAs (lncRNAs), were forming parts of the chimeric reads compared to
164 their abundance in the human genomes. These non-coding RNA chimeric sequences
165 made up 8.8% and 10.5% of the chimeric reads detected in the Caco and Calu cell
166 lines, respectively, while non-coding sequences made up only 4% of the genes
167 annotated in the human genome. This enrichment of non-coding RNA chimeric
168 sequences was tested using Fisher's exact test confirming that the trend was
169 significant (Caco cells: odds ratio=2.2, P-value=0.043; Calu cells: odds ratio=2.8,
170 P-value=0.001). When analyzed in the context of the expression level of the host
171 genes in each sample, we also observed an enrichment for highly expressed genes
172 forming parts of the chimeric sequences (Fig.2). This enrichment was confirmed
173 through the Mann-Whitney U tests showing that the trend was significant in the two
174 human cell lines (P-value < 2.2e-16 for both) and the *Chlorocebus sabaues* (green
175 monkey) cell line (P-value < 2.2e-16). These results appear to highlight two groups of

176 sequences that are forming chimeric RNAs, structural RNAs like lncRNAs, which
177 may be susceptible due to their secondary structures, and highly expressed genes,
178 which would have more RNA molecules present for template-switching events to
179 occur with. This suggests that the formation of chimeras is largely stochastic, with
180 factors like the abundance of RNAs playing a large role, but that certain RNA
181 molecules may be more susceptible to these events due to their structure.

182 **Systematic search for host-derived insertions in SARS-CoV-2 genomes**

183 We performed a survey of the GISAID SARS-CoV-2 genomes to identify insertions
184 with potential host origins. Insertions were detected based on alignments and
185 comparison to the Wuhan-Hu-1/2019 reference genome. Only insertions greater than
186 or equal to 21 nucleotides long and that were found outside of the 5' and 3'
187 untranslated regions were considered in subsequent analyses (Supplementary Table 4).
188 Of the 36 insertions that were found, 17 of them were found in multiple SARS-CoV-2
189 genomes but were not monophyletic. Upon further examination, the genomes
190 containing these insertions tended to be sequenced by the same labs around the same
191 times making it likely that these detected insertions are due to library preparation or
192 sequencing errors rather than the result of multiple independent insertion events in
193 different viral lineages. Of the 19 other insertions, 16 of them were only detected in a
194 single genome, and while many of these had plausible hits to human genes, it is
195 difficult to assess if these are true insertions or library preparation or sequencing
196 artifacts due to their limited presence.

197 The three remaining insertions were from monophyletic virus variants and were
198 further examined to determine if they had plausible homologous sequences in the
199 human genome. Two of the insertions were found to be identical to conserved
200 segments of the 28S and 18S rRNAs and were analyzed further. The remaining
201 insertion was 21 nucleotides long and was found in 6 SARS-CoV-2 genomes of the
202 Alpha B.1.1.7 lineage. These genomes were collected in early March of 2021 from
203 England, United Kingdom by two laboratories, and sequenced at the same location
204 using the same sequencing platform. The raw reads were available for two of the
205 genomes, namely England/ALDP-13C8C28/2021 (EPI_ISL_1331302) and
206 England/QEUH-13C1955/2021 (EPI_ISL_1332461), and were examined directly,
207 providing confirmation that the insertion was present and likely not an artifact.
208 Unfortunately, no plausible source for this insertion was able to be identified using a
209 BLAST search in the NCBI non-redundant nucleotide database and a collection of
210 coronavirus genomes with a cutoff E-value of 1e-2, and it was not analyzed further.

211 **28S rRNA-derived insertion in SARS-CoV-2 genomes**

212 We detected a 27-nucleotide long insertion in five SARS-CoV-2 genomes
213 (Supplementary Table 5 and Fig.3A) at position 7120 of the reference genome
214 (China/Wuhan-Hu-1/2019). The five genomes containing the 28S rRNA-derived
215 insertions were collected by different laboratories and were sequenced on different
216 sequencing platforms, making it extremely unlikely that laboratory error is
217 responsible for the presence of the insertions. The five genomes belong to a

218 monophyletic group. In this clade, there are three other variants whose assembled
219 genomes do not contain the insertion. We were able to obtain access to the raw
220 genome sequencing data of two of the three variants — USA/WA-PHL-005726/2021
221 (EPI_ISL_6259191) and USA/HI-H215617/2021 (EPI_ISL_6540096). We then did
222 further analysis on the raw sequencing to check if the insertion is indeed missing.
223 First, we generated consensus genome sequences based on the alignment of
224 sequencing reads to the SARS-CoV-2 reference genome and found the consensus
225 sequences did not contain the insertion. Next, we manually added the 28S
226 rRNA-derived insertion at position 7120 of the consensus genome and compared the
227 reads exclusively aligned to the consensus genome with the insertion and the reads
228 exclusively aligned to the consensus genome without the insertion. We found that
229 99.76% (8700/8721 for EPI_ISL_6259191) and 99.93% (1502/1503 for
230 EPI_ISL_6540096) of the exclusively-mapped reads support the presence of the
231 insertion in the genomes. The reason that the insertion is missing in the submitted
232 genomes (EPI_ISL_6259191 and EPI_ISL_6540096) is likely that the assembly was
233 generated using an insertion-unaware approach, such as reference-based consensus
234 calling. For the only one variant that was missing the insertion in the genomes, we are
235 not able to assess if it is due to failure to identify the insertion based on the consensus
236 caller or the subsequent loss of the inserted sequence.

237 By performing the BLAST search for this insertion against the human transcripts
238 (Release 109 RNAs), an exact match (E-value: 2e-06) of this insertion was found in

239 the nucleotide sequences of 28S ribosomal RNA (Fig.3B). We observed an extra three
240 overlapping bases in the pairwise alignment of SARS-CoV-2 variants containing the
241 insertion and the human 28S rRNA sequence, extending the length of identity
242 nucleotide bases from 27 nucleotides to 30 nucleotides. The identical region was
243 located at positions 4969-4998 of the human 28S rRNA (based on the structure of
244 PDB 5AJ0 Chain A2) and makes up part of the highly conserved loop 94 stem of
245 domain 7 of the rRNA molecule according to the Gorski et al.'s segmentation of
246 human 28S rRNA (26) (Fig.3B).

247 Due to the high level of sequence conservation of 28S rRNA, asserting the origin of
248 the insertion-related 30 nucleotide sequences is impossible based on sequence identity
249 alone. In the human genome (GRCh38 release 105), three 28S rRNA gene copies in
250 chromosome 21 and one copy in chromosome 12 contain the exact 30 nucleotide
251 sequences. When we searched the 30 nucleotide sequences in the LSU rRNA database
252 downloaded from SILVA (27), 98 organisms were found to contain the sequences.
253 The last common ancestor of these 98 organisms is *Euteleostomi* (bony vertebrates).
254 Given the fact that the insertion emerged from the SARS-CoV-2 variant circulating in
255 humans, the originating organism of the 28S rRNA-derived insertion is most likely
256 humans.

257 The nine amino acid insertion is located at position 1467 of the ectodomain (3Ecto) in
258 the Nsp3 protein, the only domain of this protein located on the luminal side of the
259 endoplasmic reticulum (Fig.3C). Nsp3 along with Nsp4 and Nsp6 have been shown to

260 be involved in the formation of double membrane vesicles in coronavirus infected
261 cells (28, 29). The 3Ecto domain is specifically involved in the recruitment of Nsp4
262 and has been shown to be an essential component of Nsp3 for correct
263 double-membrane vesicle formation (28). At this point, it is unclear if this insertion
264 would have had an effect on viral fitness, but given its location in the 3Ecto domain, it
265 is possible that the insertion could have an effect on the interactions between Nsp3
266 and other proteins and on the membrane rearrangement process.

267 The monophyletic group with the 28S rRNA-derived insertion belonged to the
268 AY.103 group of the delta lineage (30) (Fig.3A). The AY.103 variant was first
269 detected worldwide on January 1st, 2021 and in the USA on January 2nd, 2021. The
270 clade containing the 28S rRNA-derived insertion is defined by five nucleotide
271 mutations (T7900C, A10420T, C18646T, C25721T, and C29668T). By September
272 2021, AY.103 had become the most common delta lineage in the United States and
273 has continued to be responsible for a significant fraction of cases until the recent
274 emergence of the Omicron variant (31). The five genomes containing the 28S
275 rRNA-derived insertion were collected between October 9th and November 10th in
276 2021 from the states of Washington, Idaho, Massachusetts, and California, indicating
277 that these variants were likely being transmitted over this timeframe, but the extent to
278 which it was being spread seems to be low as Idaho was the only state where multiple
279 genomes were collected from and no genomes containing the insertion have been
280 reported since. Based on the limited spread of the viruses containing the 28S

281 rRNA-derived insertion, it is likely that the insertion might not confer phenotypic
282 advantages or is possibly disadvantageous to the virus. Nonetheless, our data show
283 that AY.103 lineages containing this insertion were viable and were transmitted for a
284 short period of time.

285 **18S rRNA-derived insertion in SARS-CoV-2 genomes**

286 A 24-nucleotide insertion was detected in two genomes at position 27492 in the
287 genome of the reference genome (China/Wuhan-Hu-1/2019) (Supplementary Table 5).
288 A sequence search against human transcripts (Release 109) was performed using
289 BLAST (32), resulting in the identification of an exact match to a 24 nucleotide
290 stretch (E-value: $2e-5$) of the 18S rRNA sequence. When aligned to the full 18S
291 rRNA sequence, it was found that the identical region extended one additional
292 nucleotide outside of the insertion region, bringing the identical stretch to 25
293 nucleotides (Fig.4A). The insertion was identical to a highly conserved region of the
294 18S rRNA (at positions 399-423 in 18S rRNA), consisting of a portion of the helix 12
295 of the 5' domain (33, 34). In the human genome alone there are five copies of the 18S
296 rRNA gene on chromosome 21 that contain identical matches for this 25 nucleotide
297 sequence. When compared to the SSU rRNA SILVA database (27), identical
298 sequences were found in the 18S sequences of 2289 organisms, which had a common
299 ancestor of *Opisthokonta* (Fungi/Metazoa group). Considering that the viral samples
300 were circulating in human populations, it is highly likely that the insertion was
301 derived from human 18S rRNA.

302 The insertion is in the SARS-CoV-2 ORF7a protein, encoding an eight amino acid
303 sequence that is located between the proline and cysteine at positions 34 and 35 in the
304 reference protein sequence (Fig.4B). The cysteine at position 35 is known to form a
305 disulfide bond with a cysteine at position 67 and is thought to help stabilize the
306 beta-sheet structure (35, 36) and the possible functions of the proline at position 34
307 are not known. The ORF7a protein has been shown to contain an
308 immunoglobulin-like ectodomain between residues 16 and 96 on the protein which is
309 thought to have a role of binding to human immune cells and modulating immune
310 response (35-37). Given the proximity of the insert to the disulfide bond forming
311 cysteine at position 34 and the size of the insert it is possible that this insert would
312 have an effect on the overall structure and immunoregulatory functions of ORF7a, but
313 without additional evidence, the effect of this insertion on the fitness of the virus
314 remains unknown.

315 The two genomes containing the 18S rRNA insertion were from the same clade in the
316 Alpha B.1.1.7 SARS-CoV-2 lineage, which was first identified in England, United
317 Kingdom in mid-December of 2020 (Fig. 4C). This variant was designated as a
318 variant of concern due to its transmissibility and large number of mutations and
319 quickly became the dominant variant in England while spreading to other countries
320 (38). The genomes containing the 18S rRNA-derived insertion, along with the other
321 four genes in the same clade, were collected in April and May of 2021 in Oregon,
322 United States. The genomes from the variants containing the insertion were collected

323 and sequenced by different labs using different sequencing platforms, making it
324 unlikely that the insertion was a sequencing or library preparation artifact. We did not
325 detect the insertion in any of the other four genomes from this clade, indicating that
326 either they do not have the insertion, they have it but it was not detected, or that the
327 insertion was only acquired in a sub-clade within this group. After May of 2021, no
328 new genomes containing this insertion were collected, indicating that the period
329 during which these lineages were circulating may have been brief. While these viral
330 variants seem to be viable and transmitted for a short period of time, the insertion
331 likely does not confer a significant advantage or may be disadvantageous for the virus
332 resulting in its limited spread.

333 **DISCUSSION**

334 Insertions in the SARS-CoV-2 genome can be introduced through multiple
335 mechanisms and have the potential to give rise to new variants with enhanced
336 infectivity, pathogenicity, and antibody escape (2, 6), but the source of these
337 insertions is often difficult to determine and has been hotly debated (5, 7). Leveraging
338 available direct RNA sequencing data and an analysis of SARS-CoV-2 genomes, we
339 have found evidence of the formation of viral-host chimeric RNA sequences and
340 described two novel human-derived genomic insertions present in circulating variants
341 of SARS-CoV-2.

342 Through our screening of direct RNA-seq data from SARS-CoV-2 infected cell lines,
343 we found that viral-host chimeric RNAs were rare but were present in approximately
344 half of the samples analyzed. The chimeric reads all contained positive-sense viral
345 RNA sequences, indicating that these chimeric sequences are not the result of the
346 integration of the viral genetic material into the host genome, which would have
347 resulted in a nearly equal mix of positive and negative sense viral sequences (17).
348 This process does appear to be stochastic in nature though, with no preference for
349 starting with host or viral sequences during chimera formation and a higher frequency
350 of chimeras being formed with highly expressed genes in the cells. The regions in the
351 RNA where these template switching events occur appears to be influenced by the
352 secondary structure of the viral RNA, possibly due to certain structures being more
353 susceptible to template switching events similar to what has been reported in previous
354 studies (2, 3). The accurate determination of the exact junction boundaries and
355 potential base-pairings were hindered by the high error rate of 14% in direct
356 RNA-sequencing data and the limited number of host-viral chimeras detected in this
357 study. The exact molecular basis for the viral-host chimera remains unclear and future
358 investigation with larger sets of error-corrected direct RNA-seq data of SARS-CoV-2
359 could be beneficial to address this question.

360 The formation of host-viral chimeric mRNAs or subgenomic RNAs could mostly be
361 transient events, not having a long-term impact on viral fitness, but the possibility of
362 human-derived insertions in the coronavirus genomes could have significant

363 implications considering the role that genomic insertions seem to have in the
364 evolution of new SARS-CoV-2 variants (5, 6). The putative 18S and 28S-derived
365 insertions were identified in circulating variants of the SARS-CoV-2, and while these
366 particular variants did not seem to spread widely, they do provide evidence that
367 human genetic material can be a source of genomic insertions in SARS-CoV-2.
368 Interestingly, rRNAs have been established to be a source of insertions in influenza
369 genomes, in some cases resulting in significantly more pathogenic viral variants (39,
370 40). It has been speculated that these recombination events often occur with host
371 rRNAs due to their abundance in the cells, the presence of recombination hotspots on
372 rRNA molecules, and the utilization of host rRNAs during viral replication (39).
373 Similar factors may play a role in the formation of these rRNA-derived insertions in
374 SARS-CoV-2, but the formation of double-membrane vesicles during SARS-CoV-2
375 would seemingly complicate this process. There may be accidental capture of host
376 RNAs inside of the double-membrane vesicles during their formation or some
377 crossover of host RNA from the cytosol, but evidence of this is lacking and warrants
378 further investigation.

379 **CONCLUSIONS**

380 Overall, our results suggest that viral-host chimeric sequences can be formed, likely
381 through stochastic RdRp template switching events. Furthermore, we have identified
382 two long insertions in SARS-CoV-2 genomes in previously circulating variants which
383 are likely derived from human ribosomal RNAs. While the source of smaller

384 insertions that are present in many SARS-CoV-2 genomes are still difficult to identify
385 due to their short lengths, these results provide evidence that bolsters the hypothesis
386 that some of them are derived from human genetic material. The mechanisms at work
387 in the formation of these chimeric RNAs and genomic insertions are still unclear but
388 warrant further study considering the potential importance of these processes in viral
389 evolution and the emergence of new variants.

390 **METHODS**

391 **Identification of host-virus chimeric reads in SARS-CoV-2 direct-RNA seq data**

392 The nanopore direct RNA-seq data from SARS-CoV-2 infected cell lines were
393 downloaded from the NCBI SRA database (Supplementary Table 1). All reads were
394 quality trimmed using NanoFilt v2.8.0 (41), to remove the first 50 nucleotides of each
395 read and require an average quality score of at least 10 over the length of the read.

396 The trimmed reads were then mapped using Minimap2 v2.23 (42) to the
397 SARS-CoV-2 reference genome (NCBI GenBank accession: NC_045512.2) (43), and
398 either a reference *Chlorocebus sabaesus* transcriptome

399 (ftp://ftp.ensembl.org/pub/release-105/fasta/chlorocebus_sabaeus/) or human

400 transcriptome (ftp://ftp.ensembl.org/pub/release-105/fasta/homo_sapiens/). The

401 mapping files were converted to the Pairwise mApping Format (PAF) using the

402 `paftools` script that is part of Minimap2 (42). Reads that mapped to both the host and

403 SARS-CoV-2 transcriptomes were extracted for analysis as potential chimeric

404 sequences. To avoid including chimeric reads that resulted from technical artifacts

405 such as those caused by misinterpretation of open-pore states by base-calling
406 softwares (19), additional quality filtering was applied to the chimeric reads. The
407 distance between the mapped regions of the virus and the host sequence on the
408 chimeric reads was required to be less than 15 nucleotides, the junction was required
409 to be formed in the middle of the genes (not within the last 50 nucleotides of the first
410 gene sequence, nor the first 50 nucleotides of the second gene sequence), and the
411 quality score within 20 bp of either side of the junction was required to be higher than
412 the 20th percentile quality score for that read.

413 **Mapping short reads to direct-RNA seq chimeric reads**

414 We collected paired-end sequencing data on five samples with corresponding direct
415 RNA-sequencing data. The short reads were first preprocessed with fastp v0.23.1 (44)
416 and then mapped to the chimeric reads from the same samples by using Minimap2
417 v2.23 (42) with options “-ax sr -w 5” to tolerate the high error rate of the Nanopore
418 direct RNA-sequencing reads (45). Read pairs spanning the junctions were detected
419 and counted with a custom script. The numbers of read pairs supporting the chimeric
420 reads are provided in Supplementary Table 2.

421 **Analysis of junction positions in relation to viral RNA secondary structure**

422 The RNA secondary structure of the SARS-CoV-2 reference genome was obtained
423 from previous studies (46, 47) and bpRNA (48) was used to assign each residue to
424 secondary structure elements. A junction site was considered in the stem if the two
425 flanking nucleotides were in the same stem. To investigate if junctions tend to happen

426 in non-stem regions, the number of junctions occurring in base-paired positions were
427 calculated and compared with a background distribution for the numbers of junctions
428 located in stems derived from a 1000-time random sampling of the same number of
429 sites along the viral RNA strand. To further examine which types of structural
430 elements are over- or under-presented at junction sites in virus-host chimeric reads
431 and in host-virus chimeric reads, one-sided Fisher's exact test was performed.

432 **Analysis of the expression level of host genes observed in chimeric reads**

433 Gene expression profiles for two SARS-CoV-2 infected Caco-2 cell line samples
434 (GSM4477888, GSM4477889), two SARS-CoV-2 infected Calu-3 cell line samples
435 (GSM4477962, GSM4477963), and three SARS-CoV-2 infected Vero-6 cell line
436 samples (GSM4916368, GSM4916369, GSM4916370) were downloaded from the
437 GEO database. The read counts of each gene were normalized by the total number of
438 reads in each sample and by the gene length (RPKM) to represent the gene expression
439 level. The background gene set was composed of all expressed protein-coding genes
440 in the cell line. To evaluate whether the expression level of the host protein-coding
441 genes in chimeric reads is significantly greater than the expression level of the
442 background gene set, a one-sided Mann-Whitney U test was performed for each
443 sample.

444 **Identification of insertions in SARS-CoV-2 genomes**

445 The SARS-CoV-2 genomes available at GISAID (<https://www.gisaid.org/>) on
446 2021-12-17 were downloaded for analysis (n=6,163,073). The sequences were then
447 processed by NextClade CLI v1.7.0 (49) which generated a multiple sequence
448 alignment against the reference genome (Wuhan-Hu-1/2019) and provided a list of
449 single nucleotide polymorphisms, insertions, and deletions associated with each
450 genome sequence. Only sequences that passed all quality controls and were assessed
451 as “good” applied by NextClade were used for further analysis (n=5,226,229).
452 Insertions greater than or equal to 21 nucleotides long and found outside of the 5’ and
453 3’ untranslated regions of the viral genomes were kept. They were searched in the
454 NCBI non-redundant nucleotide database and a collection of coronavirus genomes
455 with BLASTN (E-value \leq 1e-2) (32) to explore their possible origins.

456 **Monophyletic test**

457 To check if the insertions of interest formed a monophyletic group, all genomes that
458 contained the same insertion were analyzed using UShER: Ultrafast Sample
459 placement on Existing tRee v0.5.1 (50) against a phylogenetic tree with available
460 genomes (n=6,257,569) from GISAID, GenBank, COG-UK and CNCB generated by
461 sarscov2phylo pipeline v.13-11-20 (51). The sequences are placed within an updated
462 global subsampled SARS-CoV-2 phylogenetic tree and local subtrees are computed to
463 show more sequences with the same context of the ones being analyzed.

464 **Verification of insertions with raw sequencing data**

465 The raw genome sequencing data of USA/WA-PHL-005726/2021 (EPI_ISL_6259191)
466 and USA/HI-H215617/2021 (EPI_ISL_6540096) were analyzed to check if the
467 insertion is indeed missing. The raw sequencing reads were processed for quality
468 control using fastp v0.23.1 (44) with default parameters and mapped to the
469 SARS-CoV-2 reference genome using BWA mem v0.7.17 (52). Primer sequences in
470 reads of EPI_ISL_6259191 were soft clipped using ivar trim (parameters: -m 1 -q 0 -s
471 4 -e) and reads in amplicons with variants in primer binding sites were removed by
472 ivar removereads v1.3.1 (53). The sequencing data of EPI_ISL_6540096 were
473 preprocessed by the providing laboratory and the primers were removed. Consensus
474 genome sequences were generated based on the alignments and it was found that the
475 consensus sequences did not contain the insertion. The 28S rRNA-derived insertion
476 was manually added at position 7120 of the consensus genomes to generate consensus
477 genomes with the insertion. The alignment files were converted to FASTQ format
478 using samtools fastq command v1.14 (54) and re-aligned to the consensus genomes
479 with or without the insertion using bowtie v2.4.4 (45) (parameter: --xeq). Reads
480 exclusively aligned to the consensus genome with the insertion and exclusively
481 aligned to the consensus genome without the insertion were identified with a custom
482 script
483 ([https://github.com/ncbi/SARS2_host_derived_insertions/blob/main/verify_insertion/i](https://github.com/ncbi/SARS2_host_derived_insertions/blob/main/verify_insertion/insertion_match_reads.py)
484 [nsertion_match_reads.py](https://github.com/ncbi/SARS2_host_derived_insertions/blob/main/verify_insertion/insertion_match_reads.py)).

485 **Data availability**

486 The datasets generated in this study and scripts are available in the github repository,

487 https://github.com/ncbi/SARS2_host_derived_insertions.

488

489 **Competing interests:** The authors declare that they have no competing interests.

490 **Funding:** All authors are supported by the Intramural Research Program of the NIH,

491 National Library of Medicine.

492 **Authors' contributions:** YY was involved in the execution of the analyses,

493 interpretation of the results, and writing and revision of the manuscript. KD was

494 involved in the interpretation of the results and writing and revision of the manuscript.

495 RF was involved in the execution of the analyses and writing of the manuscript. XJ

496 was involved in the conceptualization, planning, interpretation of the results, and

497 revision of the manuscript. All authors read and approved the final manuscript.

498 **Acknowledgments:** This work utilized the computational resources of the NIH HPC

499 Biowulf cluster (<http://hpc.nih.gov>). We would like to thank Eugene V. Koonin and

500 Sofya K. Garushyants for their thoughtful comments on our manuscript and their code

501 on how to perform permutation tests and plot the results provided at

502 https://github.com/garushyants/covid_insertions_paper. We gratefully acknowledge

503 the researchers from the originating laboratories responsible for obtaining the

504 specimens and the submitting laboratories where genetic sequence data were

505 generated and shared via the GISAID Initiative, on which this research is based

506 (Supplementary Table 6). We want to particularly thank the Washington State Public

507 Health Laboratories and the State of Hawaii Laboratories Division for sharing the raw

508 sequencing data for the genomes USA/WA-PHL-005726/2021 and

509 USA/HI-H215617/2021.

510 REFERENCES

- 511 1. Gerdol M, Dishnica K, Giorgetti A. 2022. Emergence of a recurrent insertion
512 in the N-terminal domain of the SARS-CoV-2 spike glycoprotein. *Virus*
513 *Res*:198674.
- 514 2. Garushyants SK, Rogozin IB, Koonin EV. 2021. Template switching and
515 duplications in SARS-CoV-2 genomes give rise to insertion variants that merit
516 monitoring. *Communications biology* 4:1-9.
- 517 3. Chrisman BS, Paskov K, Stockham N, Tabatabaei K, Jung J-Y, Washington P,
518 Varma M, Sun MW, Maleki S, Wall DP. 2021. Indels in SARS-CoV-2 occur
519 at template-switching hotspots. *BioData Min* 14:1-16.
- 520 4. Laiton-Donato K, Franco-Muñoz C, Álvarez-Díaz DA, Ruiz-Moreno HA,
521 Usme-Ciro JA, Prada DA, Reales-González J, Corchuelo S,
522 Herrera-Sepúlveda MT, Naizaque J. 2021. Characterization of the emerging B.
523 1.621 variant of interest of SARS-CoV-2. *Infect Genet Evol* 95:105038.
- 524 5. Venkatakrishnan A, Anand P, Lenehan PJ, Suratekar R, Raghunathan B,
525 Niesen MJ, Soundararajan V. 2021. Omicron variant of SARS-CoV-2 harbors
526 a unique insertion mutation of putative viral or human genomic origin.
- 527 6. Andreano E, Piccini G, Licastro D, Casalino L, Johnson NV, Paciello I, Dal
528 Monego S, Pantano E, Manganaro N, Manenti A. 2021. SARS-CoV-2 escape
529 from a highly neutralizing COVID-19 convalescent plasma. *Proceedings of*
530 *the National Academy of Sciences* 118.
- 531 7. Peacock TP, Bauer DL, Barclay WS. Putative host origins of RNA insertions
532 in SARS-CoV-2 genomes. *Image (IN)* 556:30.4.
- 533 8. Knoop K, Kikkert M, Worm SHvd, Zevenhoven-Dobbe JC, Van Der Meer Y,
534 Koster AJ, Mommaas AM, Snijder EJ. 2008. SARS-coronavirus replication is
535 supported by a reticulovesicular network of modified endoplasmic reticulum.
536 *PLoS Biol* 6:e226.
- 537 9. Robson F, Khan KS, Le TK, Paris C, Demirbag S, Barfuss P, Rocchi P, Ng
538 W-L. 2020. Coronavirus RNA proofreading: molecular basis and therapeutic
539 targeting. *Mol Cell* 79:710-727.
- 540 10. Banner LR, Mc Lai M. 1991. Random nature of coronavirus RNA
541 recombination in the absence of selection pressure. *Virology* 185:441-445.
- 542 11. Liao C, Lai M. 1992. RNA recombination in a coronavirus: recombination
543 between viral genomic RNA and transfected RNA fragments. *J Virol*
544 66:6117-6124.
- 545 12. Yang Y, Yan W, Hall AB, Jiang X. 2021. Characterizing transcriptional
546 regulatory sequences in coronaviruses and their role in recombination. *Mol*
547 *Biol Evol* 38:1241-1248.
- 548 13. Simon-Loriere E, Holmes EC. 2011. Why do RNA viruses recombine? *Nature*
549 *Reviews Microbiology* 9:617-626.

- 550 14. Bolze A, White S, Basler T, Dei Rossi A, Roychoudhury P, Greninger AL,
551 Hayashibara K, Wyman D, Kil E, Dai H. 2022. Evidence for SARS-CoV-2
552 Delta and Omicron co-infections and recombination. medRxiv.
- 553 15. Forni D, Cagliani R, Clerici M, Sironi M. 2017. Molecular evolution of human
554 coronavirus genomes. *Trends Microbiol* 25:35-48.
- 555 16. Yan B, Chakravorty S, Mirabelli C, Wang L, Trujillo-Ochoa JL, Chauss D,
556 Kumar D, Lionakis MS, Olson MR, Wobus CE. 2021. Host-virus chimeric
557 events in SARS-CoV-2-infected cells are infrequent and artifactual. *J Virol*
558 95:e00294-21.
- 559 17. Zhang L, Richards A, Barrasa MI, Hughes SH, Young RA, Jaenisch R. 2021.
560 Reverse-transcribed SARS-CoV-2 RNA can integrate into the genome of
561 cultured human cells and can be expressed in patient-derived tissues.
562 *Proceedings of the National Academy of Sciences* 118.
- 563 18. Briggs E, Ward W, Rey S, Law D, Nelson K, Bois M, Ostrov N, Lee HH,
564 Laurent JM, Mita P. 2021. Assessment of potential SARS-CoV-2 virus
565 integration into human genome reveals no significant impact on RT-qPCR
566 COVID-19 testing. *Proceedings of the National Academy of Sciences* 118.
- 567 19. Parry R, Gifford RJ, Lytras S, Ray SC, Coin LJ. 2021. No evidence of
568 SARS-CoV-2 reverse transcription and integration as the origin of chimeric
569 transcripts in patient tissues. *Proceedings of the National Academy of*
570 *Sciences* 118.
- 571 20. Smits N, Rasmussen J, Bodea GO, Amarilla AA, Gerdes P, Sanchez-Luque FJ,
572 Ajjikuttira P, Modhiran N, Liang B, Faivre J. 2021. No evidence of human
573 genome integration of SARS-CoV-2 found by long-read DNA sequencing.
574 *Cell Rep* 36:109530.
- 575 21. Zhang L, Richards A, Barrasa MI, Hughes SH, Young RA, Jaenisch R. 2021.
576 Response to Parry et al.: Strong evidence for genomic integration of
577 SARS-CoV-2 sequences and expression in patient tissues. *Proc Natl Acad Sci*
578 *U S A* 118.
- 579 22. Shu Y, McCauley J. 2017. GISAID: Global initiative on sharing all influenza
580 data—from vision to reality. *Eurosurveillance* 22:30494.
- 581 23. Wang D, Jiang A, Feng J, Li G, Guo D, Sajid M, Wu K, Zhang Q, Ponty Y,
582 Will S. 2021. The SARS-CoV-2 subgenome landscape and its novel
583 regulatory features. *Mol Cell* 81:2135-2147. e5.
- 584 24. Wu H-Y, Brian DA. 2007. 5' -Proximal Hot Spot for an Inducible
585 Positive-to-Negative-Strand Template Switch by Coronavirus RNA-Dependent
586 RNA Polymerase. *J Virol* 81:3206-3215.
- 587 25. Kim D, Lee J-Y, Yang J-S, Kim JW, Kim VN, Chang H. 2020. The
588 architecture of SARS-CoV-2 transcriptome. *Cell* 181:914-921. e10.
- 589 26. Gorski JL, Gonzalez IL, Schmickel RD. 1987. The secondary structure of
590 human 28S rRNA: the structure and evolution of a mosaic rRNA gene. *J Mol*
591 *Evol* 24:236-251.

- 592 27. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, Peplies J,
593 Glöckner FO. 2012. The SILVA ribosomal RNA gene database project:
594 improved data processing and web-based tools. *Nucleic Acids Res*
595 41:D590-D596.
- 596 28. Hagemeyer MC, Monastyrska I, Griffith J, van der Sluijs P, Voortman J, en
597 Henegouwen PMvB, Vonk AM, Rottier PJ, Reggiori F, De Haan CA. 2014.
598 Membrane rearrangements mediated by coronavirus nonstructural proteins 3
599 and 4. *Virology* 458:125-135.
- 600 29. Lei J, Kusov Y, Hilgenfeld R. 2018. Nsp3 of coronaviruses: Structures and
601 functions of a large multi-domain protein. *Antiviral Res* 149:58-74.
- 602 30. Khare S, Gurry C, Freitas L, Schultz MB, Bach G, Diallo A, Akite N, Ho J,
603 Lee RT, Yeo W. 2021. GISAIID's Role in Pandemic Response. *China CDC*
604 *Weekly* 3:1049.
- 605 31. Hadfield J, Megill C, Bell SM, Huddleston J, Potter B, Callender C,
606 Sagulenko P, Bedford T, Neher RA. 2018. Nextstrain: real-time tracking of
607 pathogen evolution. *Bioinformatics* 34:4121-4123.
- 608 32. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K,
609 Madden TL. 2009. BLAST+: architecture and applications. *BMC*
610 *Bioinformatics* 10:1-9.
- 611 33. Gopanenko AV, Malygin AA, Karpova GG. 2015. Exploring human 40S
612 ribosomal proteins binding to the 18S rRNA fragment containing major 3'
613 -terminal domain. *Biochimica et Biophysica Acta (BBA)-Proteins and*
614 *Proteomics* 1854:101-109.
- 615 34. Granneman S, Petfalski E, Swiatkowska A, Tollervey D. 2010. Cracking pre
616 - 40S ribosomal subunit structure by systematic analyses of RNA-protein
617 cross - linking. *The EMBO journal* 29:2026-2036.
- 618 35. Cao Z, Xia H, Rajsbaum R, Xia X, Wang H, Shi P-Y. 2021. Ubiquitination of
619 SARS-CoV-2 ORF7a promotes antagonism of interferon response. *Cell Mol*
620 *Immunol* 18:746-748.
- 621 36. Zhou Z, Huang C, Zhou Z, Huang Z, Su L, Kang S, Chen X, Chen Q, He S,
622 Rong X. 2021. Structural insight reveals SARS-CoV-2 ORF7a as an
623 immunomodulating factor for human CD14+ monocytes. *IScience* 24:102187.
- 624 37. Su C-M, Wang L, Yoo D. 2021. Activation of NF- κ B and induction of
625 proinflammatory cytokine expressions mediated by ORF7a protein of
626 SARS-CoV-2. *Sci Rep* 11:1-12.
- 627 38. Volz E, Mishra S, Chand M, Barrett JC, Johnson R, Geidelberg L, Hinsley
628 WR, Laydon DJ, Dabrera G, O'Toole Á. 2021. Assessing transmissibility of
629 SARS-CoV-2 lineage B. 1.1. 7 in England. *Nature* 593:266-269.
- 630 39. Gultyaev AP, Spronken MI, Funk M, Fouchier RA, Richard M. 2021.
631 Insertions of codons encoding basic amino acids in H7 hemagglutinins of
632 influenza A viruses occur by recombination with RNA at hotspots near
633 snoRNA binding sites. *RNA* 27:123-132.

- 634 40. Khatchikian D, Orlich M, Rott R. 1989. Increased viral pathogenicity after
635 insertion of a 28S ribosomal RNA sequence into the haemagglutinin gene of
636 an influenza virus. *Nature* 340:156-157.
- 637 41. De Coster W, D'Hert S, Schultz DT, Cruts M, Van Broeckhoven C. 2018.
638 NanoPack: visualizing and processing long-read sequencing data.
639 *Bioinformatics* 34:2666-2669.
- 640 42. Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences.
641 *Bioinformatics* 34:3094-3100.
- 642 43. Wu F, Zhao S, Yu B, Chen Y-M, Wang W, Song Z-G, Hu Y, Tao Z-W, Tian
643 J-H, Pei Y-Y. 2020. A new coronavirus associated with human respiratory
644 disease in China. *Nature* 579:265-269.
- 645 44. Chen S, Zhou Y, Chen Y, Gu J. 2018. fastp: an ultra-fast all-in-one FASTQ
646 preprocessor. *Bioinformatics* 34:i884-i890.
- 647 45. Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2.
648 *Nature methods* 9:357-359.
- 649 46. Cao C, Cai Z, Xiao X, Rao J, Chen J, Hu N, Yang M, Xing X, Wang Y, Li M.
650 2021. The architecture of the SARS-CoV-2 RNA genome inside virion.
651 *Nature communications* 12:1-14.
- 652 47. Huston NC, Wan H, Strine MS, Tavares RdCA, Wilen CB, Pyle AM. 2021.
653 Comprehensive in vivo secondary structure of the SARS-CoV-2 genome
654 reveals novel regulatory motifs and mechanisms. *Mol Cell* 81:584-598. e5.
- 655 48. Danaee P, Rouches M, Wiley M, Deng D, Huang L, Hendrix D. 2018. bpRNA:
656 large-scale automated annotation and analysis of RNA secondary structure.
657 *Nucleic Acids Res* 46:5381-5394.
- 658 49. Aksamentov I, Roemer C, Hodcroft EB, Neher RA. 2021. Nextclade: clade
659 assignment, mutation calling and quality control for viral genomes. *Journal of*
660 *Open Source Software* 6:3773.
- 661 50. Turakhia Y, Thornlow B, Hinrichs AS, De Maio N, Gozashti L, Lanfear R,
662 Haussler D, Corbett-Detig R. 2021. Ultrafast Sample placement on Existing
663 tRees (USHER) enables real-time phylogenetics for the SARS-CoV-2
664 pandemic. *Nat Genet* 53:809-816.
- 665 51. Lanfear R, Mansfield R. 2020. A global phylogeny of SARS-CoV-2
666 sequences from GISAID. Zenodo doi 10.
- 667 52. Li H, Durbin R. 2009. Fast and accurate short read alignment with
668 Burrows–Wheeler transform. *Bioinformatics* 25:1754-1760.
- 669 53. Grubaugh ND, Gangavarapu K, Quick J, Matteson NL, De Jesus JG, Main BJ,
670 Tan AL, Paul LM, Brackney DE, Grewal S. 2019. An amplicon-based
671 sequencing framework for accurately measuring intrahost virus diversity using
672 PrimalSeq and iVar. *Genome Biol* 20:1-19.
- 673 54. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G,
674 Abecasis G, Durbin R. 2009. The sequence alignment/map format and
675 SAMtools. *Bioinformatics* 25:2078-2079.

677 **FIGURE LEGENDS**

678 **Fig.1 Locations of the chimeric read junction sites and permutation tests for the**
679 **number of junction sites in stems.** Diagrams show how frequently junction sites
680 occur at each position on the SARS-CoV-2 genome for (A) 5'-human-SARS2-3', (B)
681 5'-SARS2-human-3', (C) 5'-monkey-SARS2-3', and (D) 5'-SARS2-monkey-3'
682 chimeric reads. Positions are colored based on the secondary structure of the
683 SARS-CoV-2 RNA, with red lines indicating that the position is in the non-stem
684 region, while gray indicates that the position is located in the stem region. Histograms
685 following each diagram show the corresponding results of permutation tests used to
686 test if the junction sites of chimeric reads are within base-paired regions of the viral
687 RNA. Each test consists of 1000 permutations and the actual frequency of junction
688 sites occurring in the stem regions is marked with a vertical red line.

689 **Fig.2 The expression level of host genes observed in chimeric reads.** The
690 expression level of host protein-coding genes observed in chimeric reads is
691 significantly higher than the background protein-coding gene expression level based
692 on studies on (A) *Homo sapiens* Caco-2 cell line, (B) *Homo sapiens* Calu-3 cell line,
693 and (C) *Chlorocebus sabaeus* Vero-6 cell line.

694 **Fig.3 The 28S rRNA-derived insertion in SARS-CoV-2 genomes.** (A) The
695 phylogeny tree shows the genomes containing the human 28S-derived insertion. The
696 clade where the insertion was detected is highlighted with a red box and the genomes
697 with the insertion are marked with red circles at the tips. The asterisk (*) indicates that
698 the insertion should be present in the variant based on raw sequencing data. (B) The
699 insertion in SARS-CoV-2 genomes potentially originate from the host 28S rRNA
700 shown by the sequence alignment of SARS-CoV-2 reference genome (NCBI
701 accession: NC_045512.2, GISAID accession: China/Wuhan-Hu-1/2019) (pink),
702 USA/CA-CDC-FG-169171/2021 (NCBI accession: OL591909.1, GISAID accession:
703 EPI_ISL_6624703) (pink) and human 28S rRNA (chain A2 of PDB 5AJ0) (blue).
704 There are five possible alignments for mapping this insertion to the reference. Only
705 the alignment with the sequence inserted after the 3rd position of 2285th codon in
706 ORF1ab is shown. The putative insertion origin is colored in red. The numbers listed
707 above and below the alignment indicate the positions of aligned bases in the original
708 sequences. The insertion sequence (red) was mapped to the 28s rRNA (blue) in a
709 human polysome 3D structure (PDB: 5AJ0). A zoom-in view of the RNA secondary
710 structure shows that the insertion is located on the No. 94 stem of domain 7 (position:
711 4969-4998) 28S rRNA region (highlighted red). (C) Diagram shows the position of
712 the human 28S rRNA-derived insertion in the ectodomain (3Ecto) of Nsp3 protein.

713 **Fig.4 The 18S rRNA-derived insertion in SARS-CoV-2 genomes.** (A) The
714 insertion in SARS-CoV-2 genomes potentially originates from the host 18S rRNA
715 shown by the sequence alignment of SARS-CoV-2 reference genome (NCBI

716 accession: NC_045512.2, GISAID accession: China/Wuhan-Hu-1/2019) (pink),
717 USA/OR-OSPHL00675/2021 (GISAID accession: EPI_ISL_2339305) (pink) and
718 human 18S rRNA (purple). The putative insertion origin is colored in red. The
719 numbers listed above and below the alignment indicate the positions of aligned bases
720 in the original sequences. The insertion sequence (red) was mapped to the 18S rRNA
721 (purple) in a human polysome 3D structure (PDB: 5AJ0). A zoom-in view of the
722 RNA secondary structure shows that the insertion covers parts of helices 11 and 12 of
723 the 5' domain of the 18S rRNA. The location of the putative insertion sequence is
724 highlighted red. **(B)** Diagram shows the position of the human 18S-derived insertion
725 on the structure of the SARS-CoV-2 ORF7a protein (PDB: 7CI3). **(C)** The phylogeny
726 tree shows the genomes containing the human 18S-derived insertion. The clade where
727 the insertion was detected is highlighted with a red box and the genomes with the
728 insertion are marked with red circles.

729 **SUPPLEMENTAL MATERIAL**

730 **Supplementary Tables**

731 Supplementary Table 1: Direct RNA-seq data analysis. Metadata associated with all
732 30 of the analyzed directed RNA-seq samples is provided along with the number of
733 reads mapped to the SARS-CoV-2 transcriptome and the count and frequency of
734 chimeric reads in each sample. The counts of chimeric reads in the host to virus and
735 virus to host orientation are listed for each sample. The references DOI for each
736 sample is also listed.

737 Supplementary Table 2: Chimeric reads supported by spanning junction short reads.

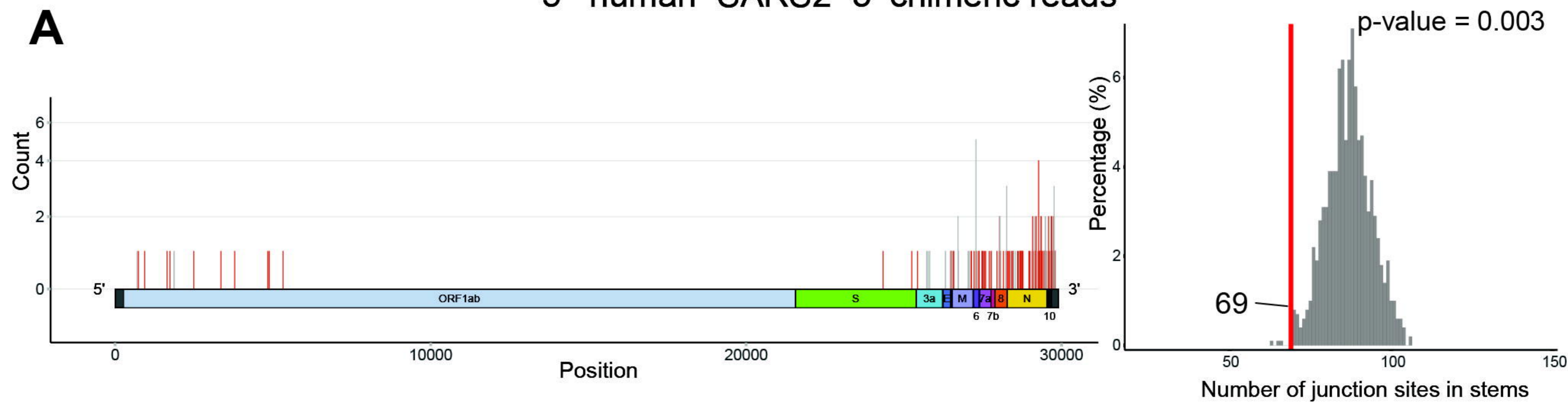
738 Supplementary Table 3: Secondary structure element enrichment analysis of
739 nucleotide sequences at junction sites.

740 Supplementary Table 4: Long insertions identified in GISAID that are not derived
741 from SARS-CoV-2. The location on the SARS-CoV-2 reference genome, insertion
742 sequence, insertion length, and what gene they are located in are provided for each of
743 the 36 detected insertions. SARS-CoV-2 genomes with the insertion, whether those
744 genomes were monophyletic, and short descriptions of putative matches are also
745 provided for each insertion.

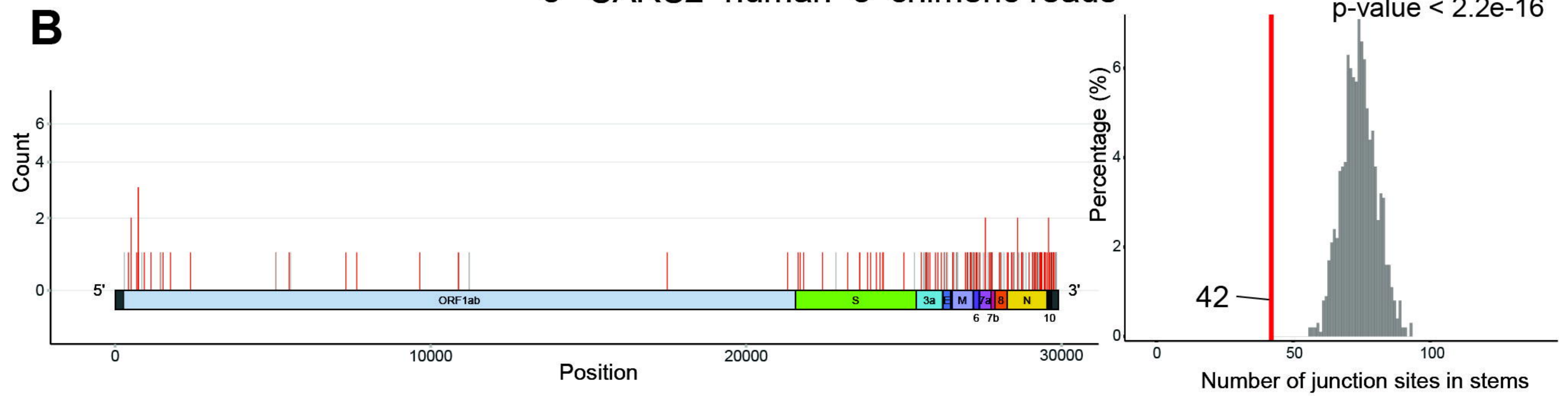
746 Supplementary Table 5: Information on genomes related to the three verified
747 insertions. Metadata associated with each of the SARS-CoV-2 genomes with putative
748 insertions that were analyzed including the variant types, collection dates, collecting
749 labs, and sequencing methods are provided.

750 Supplementary Table 6: GISAID acknowledgement table.

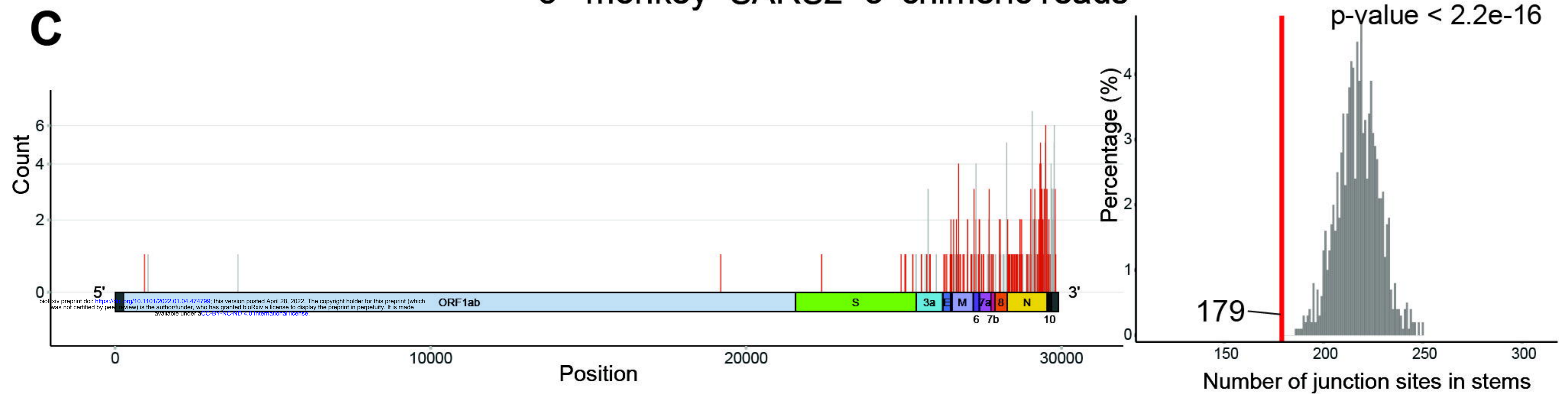
5'-human-SARS2-3' chimeric reads



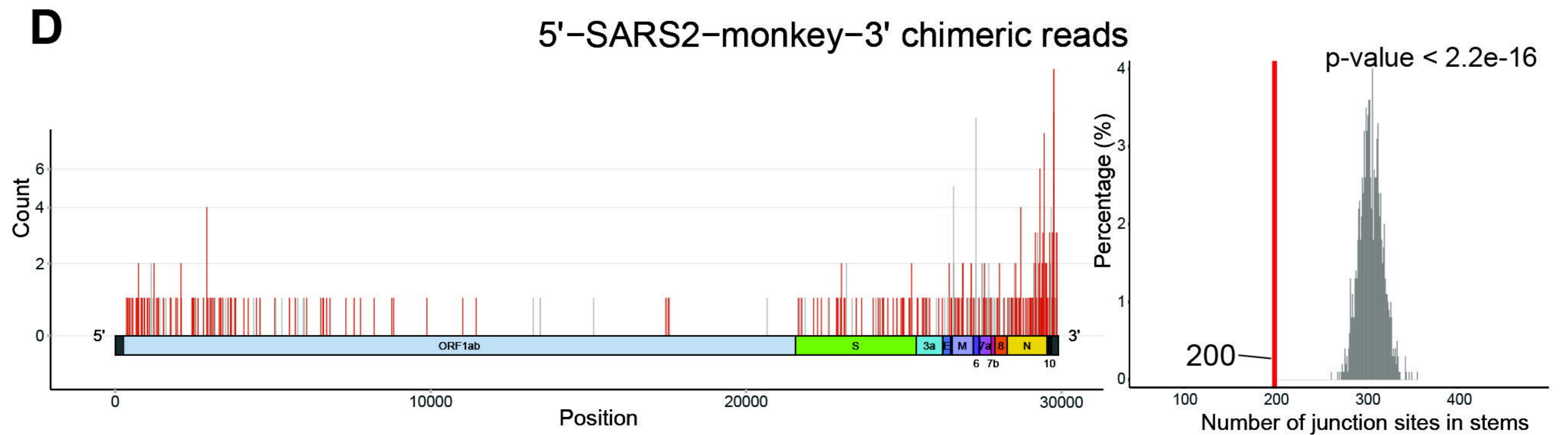
5'-SARS2-human-3' chimeric reads

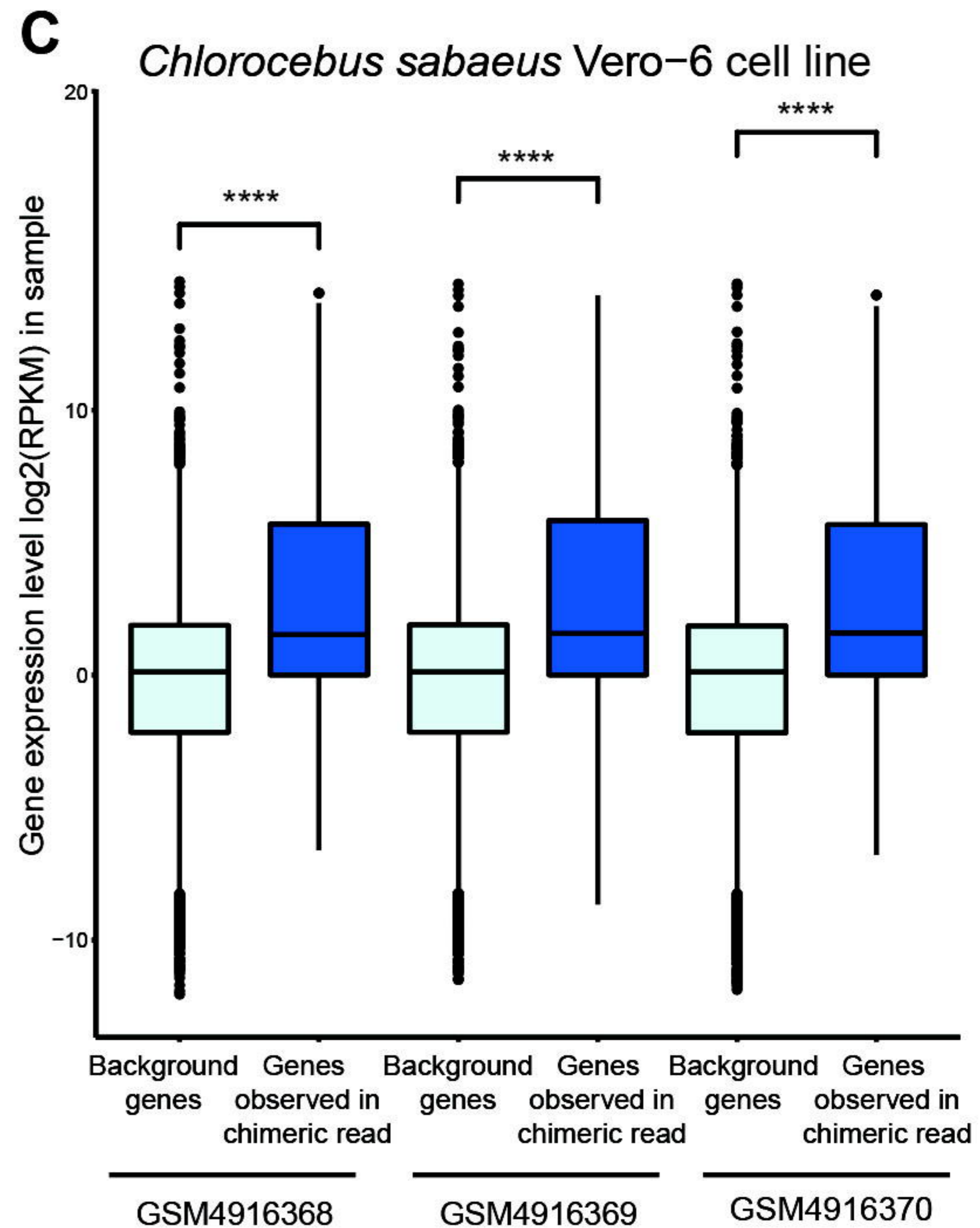
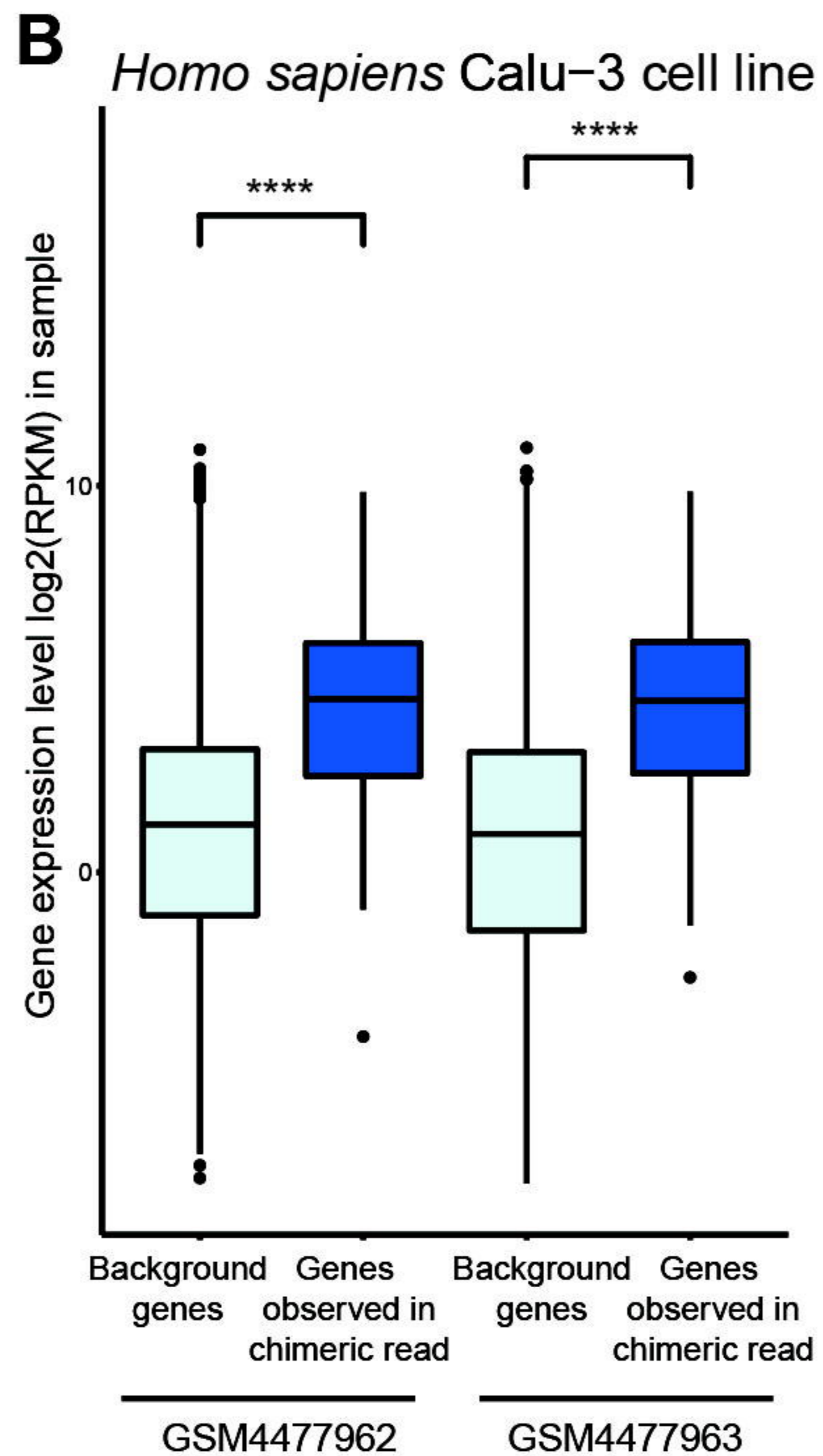
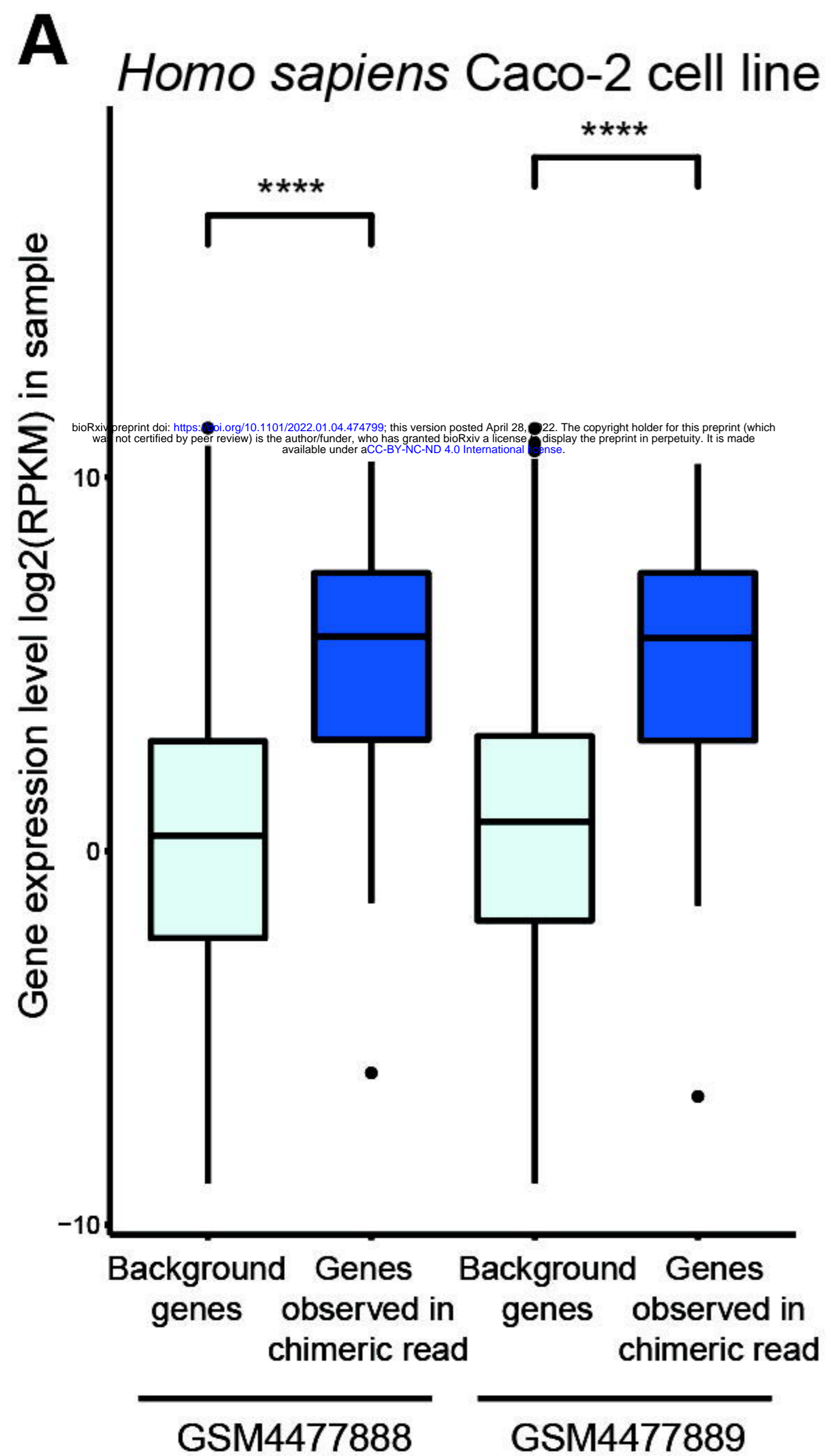


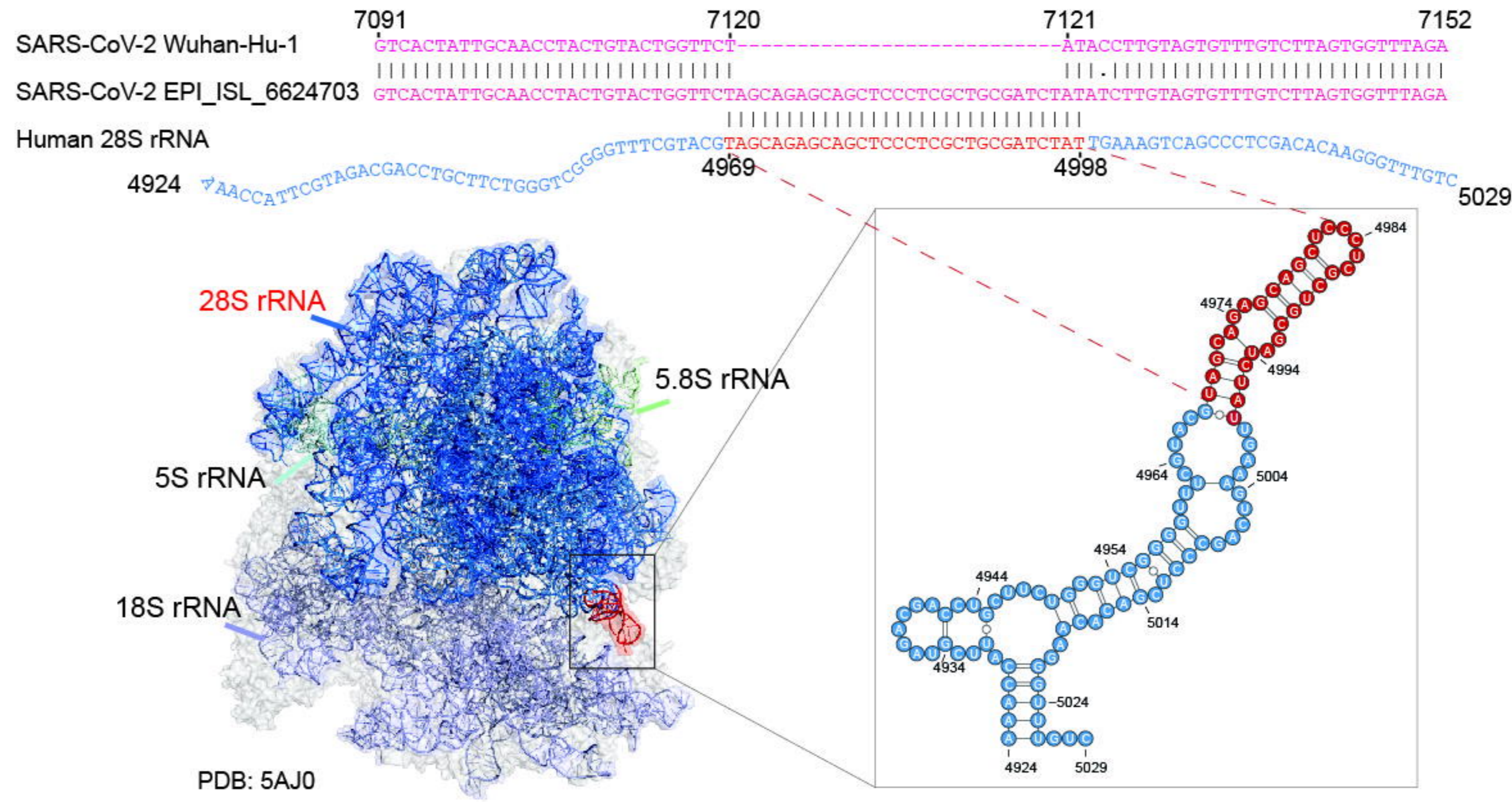
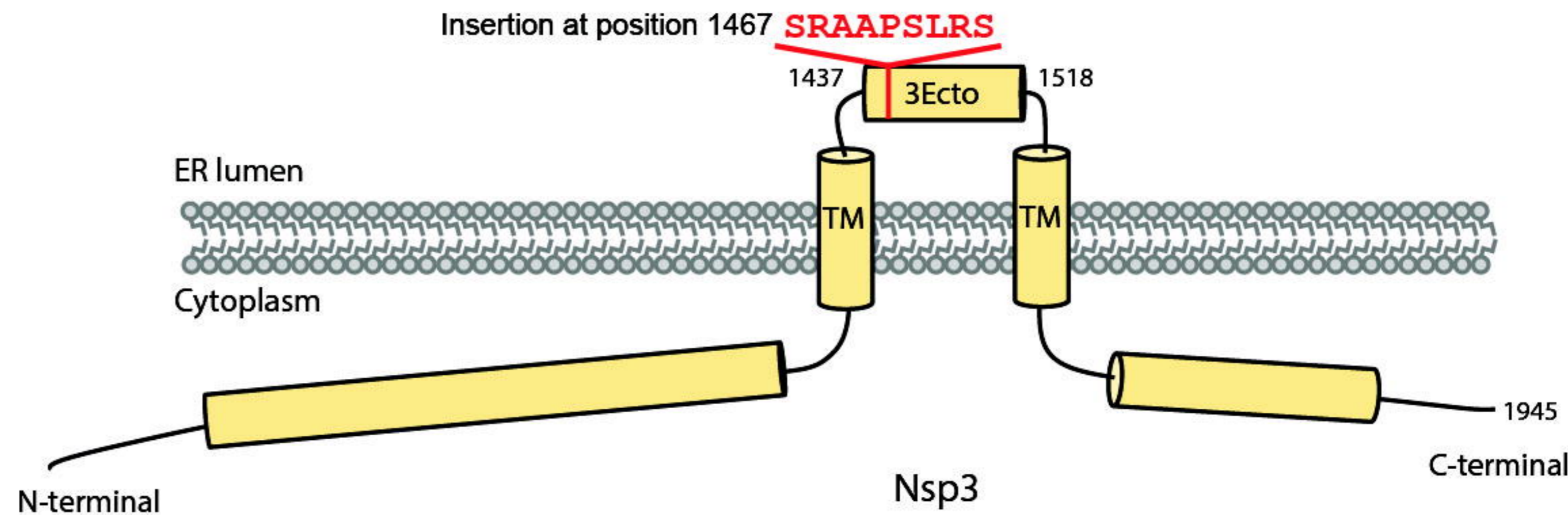
5'-monkey-SARS2-3' chimeric reads

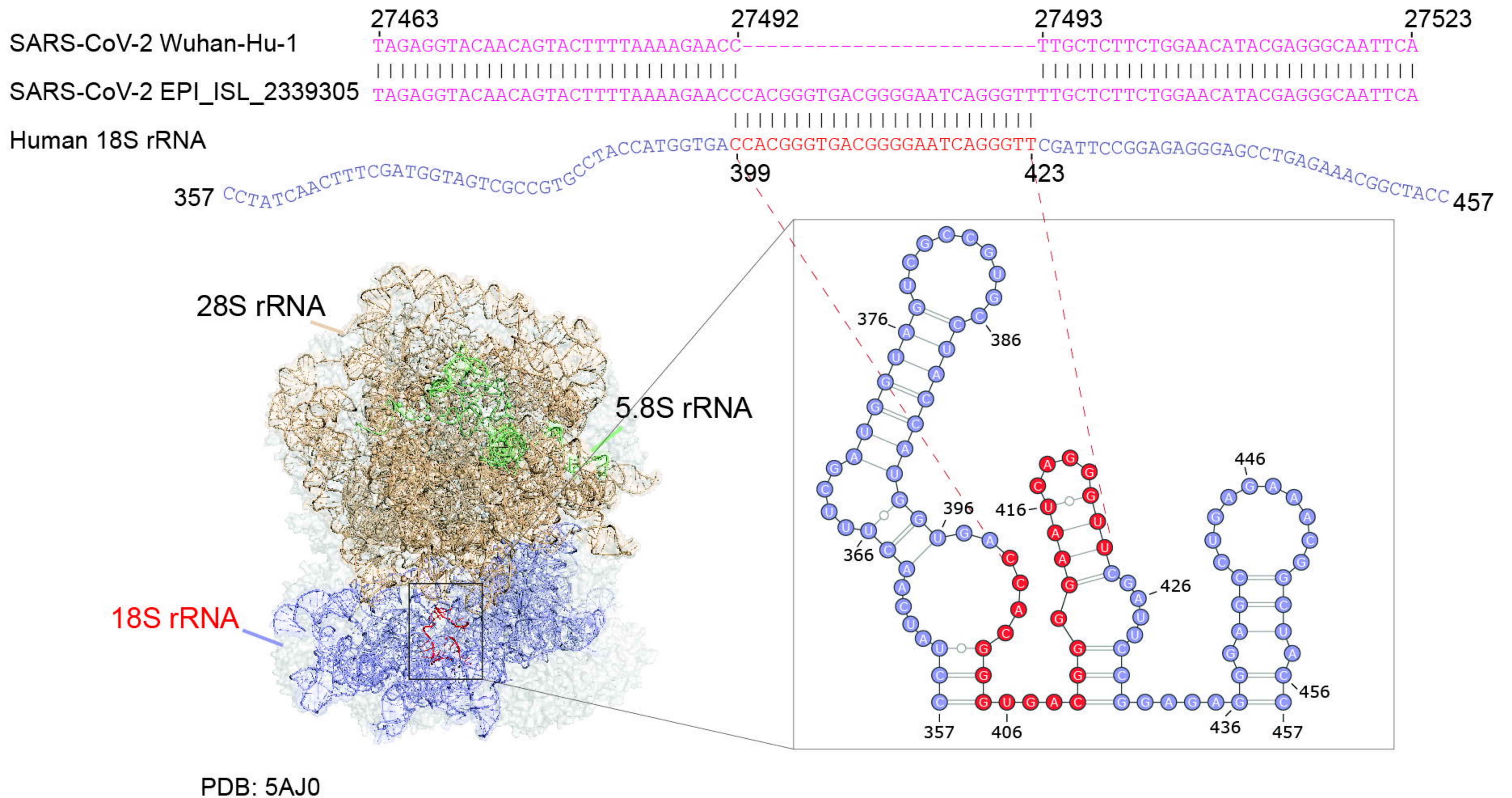
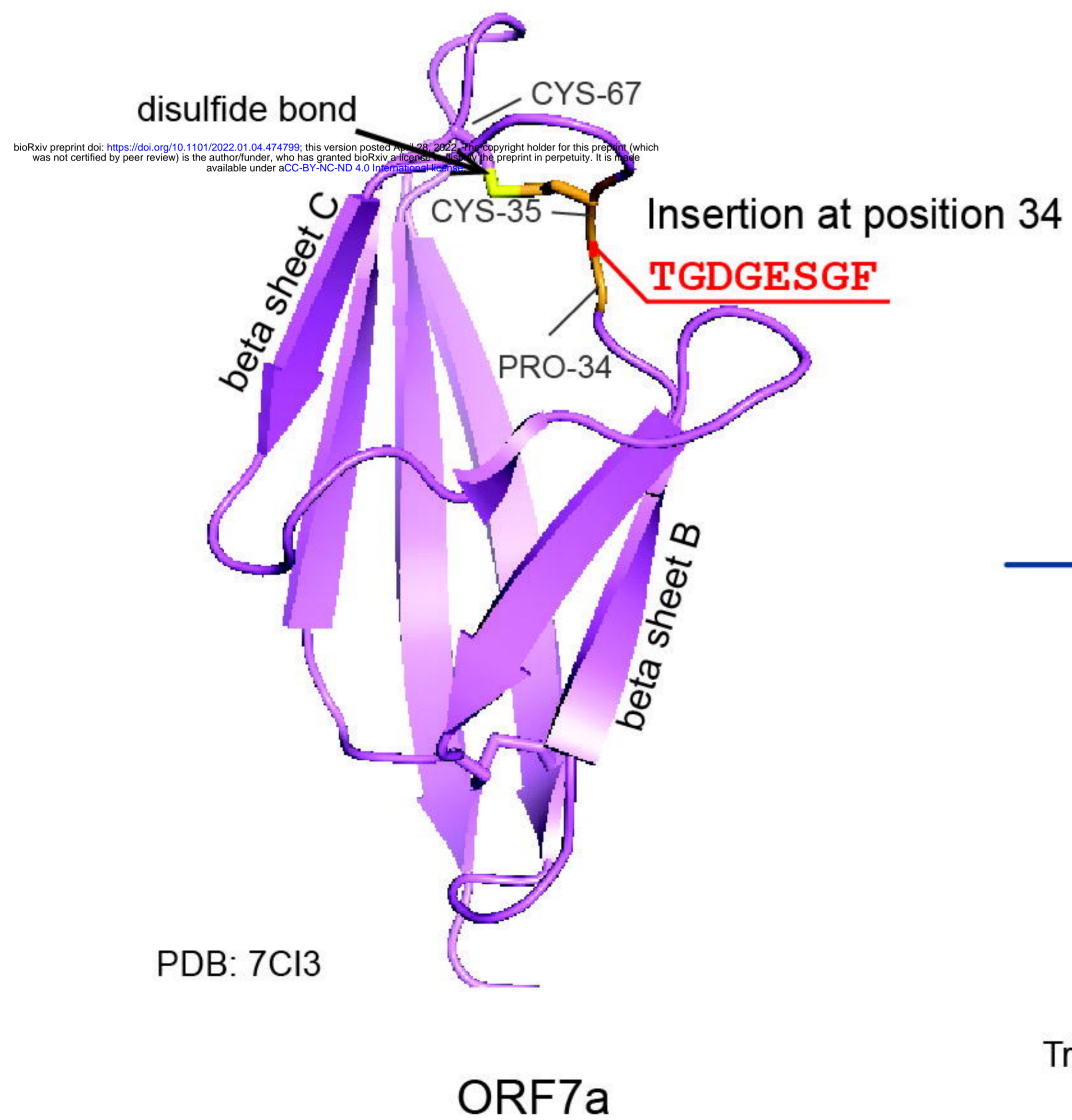
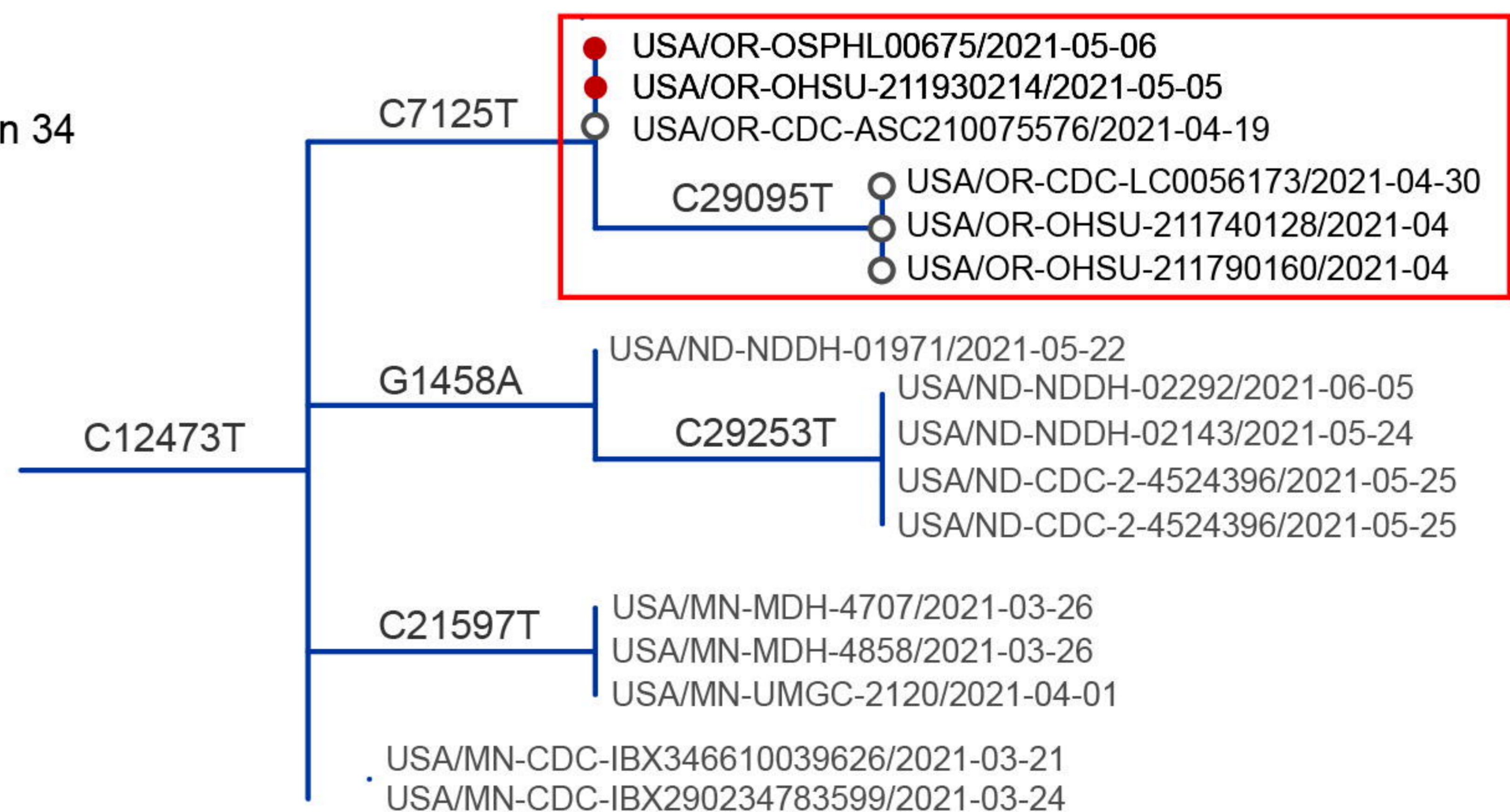


5'-SARS2-monkey-3' chimeric reads





A**B****C**

A**B****C**

Tree scale: 1