

1 **Identification of common deletions in the spike protein of SARS-CoV-2**

2 Zhe Liu<sup>\*1,2</sup>, Huanying Zheng<sup>\*2</sup>, Huifang Lin<sup>\*1,2</sup>, Mingyue Li<sup>3</sup>, Runyu Yuan<sup>1,2</sup>, Jingju Peng<sup>1,2</sup>, Qianlin

3 Xiong<sup>1,2</sup>, Jiufeng Sun<sup>1,2</sup>, Baisheng Li<sup>2</sup>, Jie Wu<sup>2</sup>, Lina Yi<sup>1,2</sup>, Xiaofang Peng<sup>1,2</sup>, Huan Zhang<sup>1,2</sup>, Wei

4 Zhang<sup>1,2</sup>, Ruben J.G. Hulswit<sup>4</sup>, Nick Loman<sup>5</sup>, Andrew Rambaut<sup>6</sup>, Changwen Ke<sup>2</sup>, Thomas A. Bowden<sup>4</sup>,

5 Oliver G Pybus<sup>7</sup>, Jing Lu<sup>1,2</sup>

6 **Affiliations:**

7 <sup>1</sup>Guangdong Provincial Institution of Public Health, Guangzhou, China;

8 <sup>2</sup>Guangdong Provincial Center for Disease Control and Prevention, Guangzhou, China;

9 <sup>3</sup>Department of Rehabilitation Medicine, The Third Affiliated Hospital, Sun Yat-sen University,

10 Guangzhou, China

11 <sup>4</sup>Division of Structural Biology, Wellcome Centre for Human Genetics, University of Oxford, Oxford,

12 UK

13 <sup>5</sup>Institute of Microbiology and Infection, University of Birmingham, UK

14 <sup>6</sup>Institute of Evolutionary Biology, University of Edinburgh, UK

15 <sup>7</sup>Department of Zoology, University of Oxford, Oxford, UK

16

17 <sup>\*</sup>Zhe Liu, Huanying Zheng, Huifang Lin contributed equally to this work. Author order was determined

18 in order of increasing seniority.

19

20 Correspondence to: Oliver G Pybus, oliver.pybus@zoo.ox.ac.uk or Jing Lu, Jimlu0331@163.com.

21 **Abstract**

22 SARS-CoV-2 is a novel coronavirus first identified in December 2019. Notable features make  
23 SARS-CoV-2 distinct from most other previously-identified Betacoronaviruses, including the receptor  
24 binding domain of SARS-CoV-2 and a unique insertion of twelve nucleotide or four amino acids  
25 (PRRA) at the S1/S2 boundary. In this study, we identified two deletion variants of SARS-CoV-2 that  
26 either directly affect the polybasic cleavage site itself (NSPRRAR) or a flanking sequence (QTQTN).  
27 These deletions were verified by multiple sequencing methods. *In vitro* results showed that the deletion  
28 of NSPRRAR likely does not affect virus replication in Vero and Vero-E6 cells, however the deletion  
29 of QTQTN may restrict late phase viral replication. The deletion of QTQTN was detected in 3 of 68  
30 clinical samples and half of 24 *in vitro* isolated viruses, whilst the deletion of NSPRRAR was identified  
31 in 3 *in vitro* isolated viruses. Our data indicate that (i) there may be distinct selection pressures on  
32 SARS-CoV-2 replication or infection *in vitro* and *in vivo*, (ii) an efficient mechanism for deleting this  
33 region from the viral genome may exist, given that the deletion variant is commonly detected after two  
34 rounds of cell passage, and (iii) the PRRA insertion, which is unique to SARS-CoV-2, is not fixed  
35 during virus replication *in vitro*. These findings provide information to aid further investigation of  
36 SARS-CoV-2 infection mechanisms and a better understanding of the NSPRRAR deletion variant  
37 observed here.

38

39 **Important notes**

40 The spike protein determines the infectivity and host range of coronaviruses. SARS-CoV-2 has two  
41 unique features in its spike protein, the receptor binding domain and an insertion of twelve nucleotides  
42 at the S1/S2 boundary resulting a furin-like cleavage site. Here, we identified two deletion variants of

43 SARS-CoV-2 that either directly affect the furin-like cleavage site itself (NSPRRAR) or a flanking  
44 sequence (QTQTN) and investigated these deletions in cell isolates and clinical samples. The absence  
45 of the polybasic cleavage site in SARS-CoV-2 did not affect virus replication in Vero or Vero-E6 cells.  
46 Our data indicate the PRRAR and its flanking sites are not fixed *in vitro*, thus there appears to be  
47 distinct selection pressures on SARS-CoV-2 sequences *in vitro* and *in vivo*. Further investigation of the  
48 mechanism of generating these deletion variants and their infectivity in different animal models would  
49 improve our understanding of the origin and evolution of this virus.  
50

51 **Introduction**

52 SARS-CoV-2 is a novel coronavirus that was first identified at the end of December 2019 (1) and  
53 responsible for the global pandemic of COVID-19(2). Unlike the two other zoonotic coronaviruses,  
54 SARS-CoV-1 and MERS-CoV(3), the evolutionary history of SARS-CoV-2 is largely unknown. A  
55 recent analysis of genetic information and the spike (S) protein structure(4, 5) highlights two notable  
56 features of the SARS-CoV-2 genome. First, the receptor binding domain (RBD) of SARS-CoV-2 is  
57 distinct from the most closely-related virus (RaTG13) of bat origin and more closely related to  
58 pangolin coronaviruses(6, 7). The spike protein of SARS-CoV-2 is demonstrated to have a high affinity  
59 for the human ACE2 receptor molecule(4). Second, a unique insertion of 12 nucleotides (encoding four  
60 amino acids, PRRAR) at the S1/S2 boundary(8) leading to a predictively solvent-exposed PRRAR/SV  
61 sequence, which corresponds to a canonical furin-like cleavage site(9, 10).

62

63 With respect to the first feature, an RBD identified in a SARS-like virus from a pangolin suggests that  
64 an RBD similar to that of SARS-CoV-2 may already exist in mammalian host(s) prior to its  
65 introduction into humans(7). The question remaining is the history and function of the insertion at the  
66 S1/S2 boundary, which is unique to SARS-CoV-2. By sequencing the whole genome of SARS-CoV-2  
67 from cell isolates and clinical samples, we identified two deletion variants that directly affect the furin  
68 cleavage site itself (NSPRRAR) or a flanking sequence (QTQTN). We screen these two deletions in  
69 cell-isolated strains and clinical samples. To explore the potential effect of these deletions, these two  
70 deletion variants were isolated and their replication kinetics were investigated in both Vero and  
71 Vero-E6 cells.

72 **Results**

73 **Identification of deletions in SARS-CoV-2 spike protein**

74 The first COVID-19 clinical case (Sample 014, Table1) in Guangdong was reported on 19<sup>th</sup> January,  
75 with illness onset on 1<sup>st</sup> January(11). A BALF (bronchoalveolar lavage fluid) sample from this patient  
76 was collected and inoculated on Vero-E6 cells. A cell-isolated viral strain was obtained after three  
77 rounds of passage. Multiple sequencing methods were used for whole genome sequencing and the  
78 validation of variants (Figure1 A, Table1), including multiplex-PCR with Miseq platform (PE150),  
79 direct CDNA sequencing in Nanopore platform and Sanger sequencing (See Materials and Methods for  
80 detail). After mapping to the SARS-CoV-2 reference genome (MN908947.3), we found that there were  
81 two variants in the cell-isolated viral strain with deletions at (1) 23585–23599 (Var1), flanking the  
82 polybasic cleavage site, resulting in a QTQTN deletion in the spike protein (one amino acid before the  
83 polybasic cleavage site) and (2) 23597–23617 (Var2), resulting in a NSPRRAR deletion that includes  
84 the polybasic cleavage site (Figure 1A). To exclude the possibility that these findings were caused by  
85 errors in PCR amplification, both of the deletion variants were verified through direct cDNA  
86 sequencing on the ONT nanopore platform. Sanger sequencing with specific primers also identified  
87 heterozygous peaks with distinct double peaks starting at the position 23585 and triple peaks after that,  
88 highlighting the existence of multiple variants caused by the above two deletions (Figure 1B). To  
89 investigate the dynamics of these deletion variants, we performed nanopore sequencing on the 014  
90 viral strain, isolated at different rounds of passage from the Vero-E6 cell culture (Figure 1C). High  
91 frequencies of the deletion variant Var1 were observed after the first passage and high frequencies of the  
92 deletion variant Var2 were observed after the 4th passage, at which point the frequency of Var1 and  
93 Var2 reached around 50%. The percentages of these two deletion variants were steady in the following

94 passages.

95

#### 96 **The deletion is commonly identified in cell isolated strains**

97 To investigate whether the deletions described above were random mutations that occasionally arise in  
98 a strain, or whether they commonly occur after cell passages, we performed whole genome sequencing  
99 on 23 other SARS-CoV-2 strains collected after two rounds of cell passage in Vero-E6 or Vero cells  
100 (Table 1). The corresponding original samples for these strains were collected between 19<sup>th</sup> January and  
101 28<sup>th</sup> February 2020. In addition to the 014 strain mentioned above, 10 out of 18 Vero-E6 isolated strains  
102 and 1 out of 5 Vero isolated strains displayed the Var1 deletion variant (>10% of sequencing reads;  
103 Figure 1D). Additionally, in two Vero-E6 isolated strains (619 and 4276), Var2 was detected, and this  
104 variant has been independently identified by another group almost at the same time, using direct RNA  
105 sequencing method(12). To find out whether these deletions were restricted to a specific genetic lineage,  
106 we next investigated the phylogenetic relationship of these viral strains. As shown in Figure 1D, the  
107 strains with a relatively higher ratio of this deletion were dispersed in the phylogenetic tree, that  
108 suggesting the deletion mutations did not arise through shared ancestry and were not restricted to a  
109 specific genetic lineage of SARS-CoV-2 viruses.

110

#### 111 **Replication kinetics of the deletion variants**

112 To evaluate the effect of these deletions on virus replication, we performed plaque assays and picked  
113 individual clones for different variants. Single plaques for Var1 and Var2 were obtained and confirmed  
114 by whole genome sequencing (014-Var1, 014-Var2; Table 1). However, the 014 strain without these  
115 deletions could not be successfully selected from plaques, possibly due to the replication advantage of

116 the deletion variants in cell culture. We investigated the replication kinetics of 014-Var1 and 014-Var2  
117 in Vero-E6 and Vero cells. The strain 029/E6 was used as a reference, which has no deletion mutations  
118 and only one amino acid difference from strain 014 on the spike protein (H47Y). The viral replication  
119 kinetics were assessed by detecting the intracellular viral loads at 1, 3, 6, 9, 12 and 24 hours post  
120 inoculation (Figure 2). As shown in Figure 2A, the 014-Var1 and 014-Var2 exhibit similar replication  
121 dynamics to the 029 strain in Vero-E6 cells. In contrast, the deletion of 23583–23599 in SARS-CoV-2  
122 (Var1) significantly diminishes cellular viral load at 24 hours post-inoculation in Vero cells (Figure 2B)  
123 and to a lesser extent in Vero-E6 cells (Figure 2A). This is the possible reason that 014-Var1 was  
124 observed less often in Vero cells than in Vero-E6 cells (Figure 1D).

125

#### 126 **Screening for deletion variants in original clinical samples**

127 To identify whether these deletions also occurred in the original clinical samples, we screened  
128 high-throughput sequencing data from 149 clinical samples, which were collected between 6<sup>th</sup> February  
129 and 20<sup>th</sup> March in Guangdong, China. There were 68 SARS-CoV-2 genomes, with an average  
130 sequencing depth  $\geq 20$  at the sites neighboring 23585. As shown in Table 2, variants with the QTQTN  
131 (Var1) were found in 3 (4%) of clinical samples, with the ratio of deletion variant in total reads ranging  
132 from 8.8–32.8%, indicating that this deletion also occurs in *in vivo* infections. Notably, two out of the  
133 three patients from which these samples were derived displayed mild symptoms and recurrence of  
134 SARS-CoV-2 infection after being discharged from hospital. The sequenced samples were collected at  
135 4 days and 17 days after discharge, respectively. The third case (20SF5645) was an asymptomatic  
136 infection case. To date, there are no genome sequences deposited in public databases containing these  
137 two deletions. While the described Var1 deletion variant was only detected in clinical samples after

138 deep sequencing, such variants may be underrepresented in databases due to the low frequency and  
139 consequent elimination upon consensus sequence generation.

140

#### 141 **Discussion**

142 The spike protein of coronaviruses plays an important role in viral infectivity, transmissibility and  
143 antigenicity. Therefore, the genetic character of the spike protein in SARS-CoV-2 may shed light on its  
144 origin and evolution(7, 8). For SARS-CoV-1, positive selection was identified in the spike coding  
145 sequence(13) and deletions in *ORF8*(14) during the early, but not late, stage of the epidemic,  
146 suggesting that SARS-CoV-1 may have been sub-optimal in the human population during the early  
147 epidemic stage after it was first transmitted from an intermediate animal host, and underwent further  
148 adaptation. SARS-CoV-2, however, has presented high infectivity and efficient transmission capability  
149 since its identification(1) suggesting the polybasic cleavage site is an important component of the virus'  
150 fitness within the human population. Genetic changes related to viral fitness of SARS-CoV-2 require  
151 further epidemiological investigation and functional analysis.

152

153 Here, we use different sequencing methods to identify and verify two deletion variants either directly  
154 affecting the polybasic cleavage site (Var1) or a site immediately upstream of it (Var2). The QTQTN  
155 deletion variant (Var1) was detected in 3 out of 68 clinical samples and half of the 24 *in vitro* isolated  
156 viral strains tested in this study. The cellular replication kinetic data suggests the deletion of the  
157 polybasic cleavage site does not affect SARS-CoV-2 replication in Vero and Vero-E6 cells, whilst the  
158 QTQTN deletion may restrict virus replication in Vero cells at the late phase. These data indicate that (i)  
159 the deletions of QTQTN and the polybasic cleavage site are likely under strong purifying selection *in*

160 *vivo*, since the deletion is rarely identified in clinical samples, (ii) there may be an efficient mechanism  
161 for generating these deletions, given that the QTQTN deletion (Var1) is commonly detected after two  
162 rounds of cell passage and (iii) the PRRA insertion, which distinguishes SARS-CoV-2 from other  
163 SARS-like viruses, is not fixed *in vitro*, because the NSPRRAR deletion variant (Var2) is observed in 3  
164 out of 24 Vero-E6 isolated strains, but does appear to be subject to purifying selection *in vivo*.

165

166 Given that these residues are located in solvent-accessible loops of the spike protein, combined with  
167 the observation that they are either partially (QTQTN) or completely (NSPRRAR) unresolved in  
168 recently reported SARS-CoV-2 S cryoEM structures(4, 5) (Figure 3), it seems likely that this region is  
169 structurally tolerant to deletions. Whilst the deletion of the furin site, as observed in Var2, would result  
170 in a loss of susceptibility to furin cleavage at this site, the effect of Var1 on furin cleavage is less  
171 evident. However, it is likely that these overlapping deletion variants have arisen through the same  
172 selective pressure and are therefore both likely to compromise furin-mediated cleavage at this position  
173 in the S protein, albeit possibly to different extents. Furthermore, it is possible that the presence of a  
174 conserved cathepsin L site 10 residues downstream of the polybasic cleavage site may provide  
175 functional tolerance(15) to any reduction in proteolytic cleavage efficiency that may arise from changes  
176 in this region (Figure 1A). Consistent with the modeling analysis, the replication dynamics in Vero and  
177 Vero-E6 cells also indicate that polybasic cleavage site deletion (Var2) does not affect virus replication  
178 *in vitro*.

179

180 Notably, a recently reported SARS-like strain, RmYN02, which is phylogenetically related to  
181 SARS-CoV-2, also has a possible deletion at the QTQT site(16). This raises another possible scenario,

182 which is that some SARS-CoV-2-like viruses in animals may not have had QTQTN in their spike  
183 protein. The origin of polybasic cleavage site (PRRA) is important to understanding the evolution  
184 history and tracing the potential animal reservoir(s) of SARS-CoV-2. Here, the different deletion  
185 frequencies observed *in vitro* and *in vivo* have provide clues that will aid further investigation of this  
186 evolutionary tale. The absence of NSPRRA in isolated SARS-CoV-2 strains could be used to further  
187 investigate its infectivity in different potential intermediate animal hosts and resolve the origin of this  
188 feature of the SARS-CoV-2 genome. In addition, the different selective pressure observed on NSPRRA  
189 region of SARS-CoV-2 *in vivo* and *in vitro* highlight the NSPRRA deletion variant generated in this  
190 study as a promising vaccine candidate in the future.  
191

192 **Materials and Methods**

193 **Ethics**

194 This study was approved by ethics committee of the Center for Disease Control and Prevention of  
195 Guangdong Province. Written consent was obtained from patients or their guardian(s) when clinical  
196 samples were collected. Patients were informed about the surveillance before providing written consent,  
197 and sequence data were analyzed anonymously.

198

199 **Viral isolation**

200 Vero E6 or Vero cells were used for SARS-CoV-2 virus isolation and passage. The cells were inoculated  
201 with 100 µl processed patient sample. Cytopathic effect (CPE) were observed daily. If there was no CPE  
202 observed, cell lysis was collected by centrifugation after three repeated freeze-thaw and 100 µl  
203 supernatant were used for the second round of passage.

204

205 **Genetic sequencing and sequence analysis**

206 The deletion variants of SARS-CoV-2 were confirmed by different approaches as previously  
207 described(17) (i) using version 1 of the ARTIC COVID-19 multiplex PCR primers  
208 (<https://artic.network/ncov-2019>), followed by sequencing on a Miseq PE150 or an ONT  
209 MinION, (ii) CDNA directly sequencing on an ONT MinION and (iii) sanger sequencing by  
210 using the nCoV-2019\_78\_LEFT and nCoV-2019\_78\_RIGTH primers from the ARTIC  
211 COVID-19 multiplex PCR primers set. The amplification products targeting the 23444-23823  
212 fragment of viral genome (numbered according to MN908947.3).

213

214 For metatranscriptomics, total RNAs were extracted from different types of samples by using  
215 QIAamp Viral RNA Mini Kit, followed by DNase treatment and purification with TURBO

216 DNase and Agencourt RNAClean XP beads. Libraries were prepared using the SMARTer  
217 Stranded Total RNA-Seq Kit v2 (according to the manufacturer's protocol starting with 10 ng  
218 total RNA. Sequencing of metatranscriptome libraries was conducted on the Illumina Miseq  
219 PE 150 platform. For the multiplex PCR approach, we followed the general method of  
220 multiplex PCR as described in (<https://artic.network/ncov-2019>)(18). Briefly, multiplex PCR was  
221 performed with two pooled primer mixtures and cDNA reverse-transcribed with random primers was  
222 used as a template. After 25-35 rounds of amplification, PCR products were collected and quantified,  
223 followed by sequencing on Illumina Miseq PE 150 platform or MinION sequencing device.  
224 Assembly of the Illumina raw data was performed using Geneious v11.0.3  
225 (<https://www.geneious.com>). Assembly of the nanopore raw data was performed using the ARTIC  
226 bioinformatic pipeline for COVID-19 with minimap2(19) and medaka  
227 (<https://github.com/nanoporetech/medaka>) for consensus sequence generation. Variant sites were called  
228 by using iVar(20) with depth  $\geq 20$  as a threshold. For direct cDNA sequencing, we followed the  
229 Nanopore Direct cDNA sequencing protocol (SQK-DCS109). Briefly, 100ng viral RNA were reverse  
230 transcribed using SuperScript™ IV First-Strand Synthesis System (Invitrogen, USA) followed by  
231 RNA chain digestion and second strand synthesis. A total of 20ng cDNA libraries were loaded to  
232 FLO-MIN106 flow cell. Generated sequences were mapped to MN908947.3 reference sequence using  
233 minimap2. The ML phylogeny for 24 viral strains genomes was estimated with PhyML(21)  
234 using the HKY+ $\Gamma_4$  substitution model(22) with gamma-distributed rate variation(23).

235

236 **Viral kinetics analysis**

237 The individual clones of deletion variants were selected by using a plaque assay. The isolated 014  
238 strains were serially-diluted and used to inoculate the monolayer of Vero-E6 cells. When CPE were  
239 observed, the cell monolayers were scraped with the back of a pipette tip. Virus lysate was used for  
240 genetic sequencing and viral strain amplification. To assess the kinetic of virus replication, different  
241 viral strains were first filtered and inoculated with Vero-E6 and Vero cells at MOI 0.5. Time was set as  
242 zero when cells were incubated with viruses. After 1 hour adsorption, the culture media were removed  
243 and cells were washed twice with PBS to remove unattached virus. Cells were lysed at different time  
244 post inoculation and total RNA was extracted by using RNeasy mini kit (QIAGEN, Germany). Cellular  
245 viral loads were calculated by using SARS-CoV-2 RT-PCR kit (DAAN GENE, Guangzhou, China) and  
246 GAPDH (glyceraldehyde-3-phosphate dehydrogenase) gene was parallelly quantified as an  
247 endogenous control.

248

#### 249 **Data Availability**

250 Metagenomic sequencing, multiplex PCR sequencing and cDNA direct sequencing data after mapping  
251 to SARS-COV-2 reference genome (MN908947.3) have been deposited in the Genome Sequence  
252 Archive(24) in BIG Data Center(25), Beijing Institute of Genomics (BIG), Chinese Academy of  
253 Sciences, under project accession numbers CRA002500, publicly accessible at  
254 <https://bigd.big.ac.cn/gsa>. The sample information and corresponding accession number for each  
255 sample are listed in the Table 1.

256

#### 257 **Acknowledgments**

258 This work was supported by grants from Guangdong Provincial Novel Coronavirus Scientific and

259 Technological Project (2020111107001), Science and Technology Planning Project of Guangdong  
260 (2018B020207006). The Wellcome Centre for Human Genetics is supported by Wellcome Centre grant  
261 203141/Z/16/Z.  
262 Conflict of interest: None declared.

263 **Reference**

- 264 1. Wu F, Zhao S, Yu B, Chen Y-M, Wang W, Song Z-G, Hu Y, Tao Z-W, Tian J-H, Pei Y-Y,  
265 Yuan M-L, Zhang Y-L, Dai F-H, Liu Y, Wang Q-M, Zheng J-J, Xu L, Holmes EC, Zhang  
266 Y-Z. 2020. A new coronavirus associated with human respiratory disease in China. *Nature*  
267 1–8.
- 268 2. WHO (2020). Coronavirus disease (COVID-2019) situation reports.
- 269 3. Cui J, Li F, Shi Z-L. 2019. Origin and evolution of pathogenic coronaviruses. *Nat Rev*  
270 *Microbiol* 17:181–192.
- 271 4. Walls AC, Park Y-J, Tortorici MA, Wall A, McGuire AT, Veesler D. 2020. Structure,  
272 Function, and Antigenicity of the SARS-CoV-2 Spike Glycoprotein. *Cell*.
- 273 5. Li F. 2005. Structure of SARS Coronavirus Spike Receptor-Binding Domain Complexed  
274 with Receptor. *Science* 309:1864–1868.
- 275 6. Xiao K, Zhai J, Feng Y, Zhou N, Zhang X, Zou J-J, Li N, Guo Y, Li X, Shen X, Zhang Z,  
276 Shu F, Huang W, Li Y, Zhang Z, Chen R-A, Wu Y-J, Peng S-M, Huang M, Xie W-J, Cai  
277 Q-H, Hou F-H, Chen W, Xiao L, Shen Y. 2020. Isolation of SARS-CoV-2-related  
278 coronavirus from Malayan pangolins. *Nature* 1–4.
- 279 7. Lam TT-Y, Shum MH-H, Zhu H-C, Tong Y-G, Ni X-B, Liao Y-S, Wei W, Cheung WY-M, Li  
280 W-J, Li L-F, Leung GM, Holmes EC, Hu Y-L, Guan Y. 2020. Identifying SARS-CoV-2  
281 related coronaviruses in Malayan pangolins. *Nature*.

- 282 8. Andersen KG, Rambaut A, Lipkin WI, Holmes EC, Garry RF. 2020. The proximal origin  
283 of SARS-CoV-2. *Nat Med*.
- 284 9. Coutard B, Valle C, de Lamballerie X, Canard B, Seidah NG, Decroly E. 2020. The spike  
285 glycoprotein of the new coronavirus 2019-nCoV contains a furin-like cleavage site absent in  
286 CoV of the same clade. *Antiviral Research* 176:104742.
- 287 10. Izaguirre G. 2019. The Proteolytic Regulation of Virus Cell Entry by Furin and Other  
288 Proprotein Convertases. *Viruses* 11:837.
- 289 11. Kang M, Wu J, Ma W, He J, Lu J, Liu T, Li B, Mei S, Ruan F, Lin L, Zou L, Ke C, Zhong  
290 H, Zhang Y, Chen X, Liu Z, Zhu Q, Xiao J, Yu J, Hu J, Zeng W, Li X, Liao Y, Tang X, Xiao  
291 S, Wang Y, Song Y, Zhuang X, Liang L, Zeng S, He G, Lin P, Deng H, Song T. 2020.  
292 Evidence and characteristics of human-to-human transmission of SARS-CoV-2. *medRxiv*  
293 2020.02.03.20019141.
- 294 12. Davidson AD, Williamson MK, Lewis S, Shoemark D, Carroll MW, Heesom K, Zambon M,  
295 Ellis J, Lewis PA, Hiscox JA, Matthews DA. 2020. Characterisation of the transcriptome  
296 and proteome of SARS-CoV-2 using direct RNA sequencing and tandem mass  
297 spectrometry reveals evidence for a cell passage induced in-frame deletion in the spike  
298 glycoprotein that removes the furin-like cleavage site. *bioRxiv* 2020.03.22.002204.
- 299 13. The Chinese SARS Molecular Epidemiology Consortium. 2004. Molecular Evolution of the  
300 SARS Coronavirus During the Course of the SARS Epidemic in China. *Science* 303:1666–  
301 1669.

- 302 14. Muth D, Corman VM, Roth H, Binger T, Dijkman R, Gottula LT, Gloza-Rausch F, Balboni  
303 A, Battilani M, Rihtarič D, Toplak I, Ameneiros RS, Pfeifer A, Thiel V, Drexler JF, Müller  
304 MA, Drosten C. 2018. Attenuation of replication by a 29 nucleotide deletion in  
305 SARS-coronavirus acquired during the early stages of human-to-human transmission.  
306 *Scientific Reports* 8:1–11.
- 307 15. Hoffmann M, Kleine-Weber H, Schroeder S, Krüger N, Herrler T, Erichsen S, Schiergens  
308 TS, Herrler G, Wu N-H, Nitsche A, Müller MA, Drosten C, Pöhlmann S. 2020.  
309 SARS-CoV-2 Cell Entry Depends on ACE2 and TMPRSS2 and Is Blocked by a Clinically  
310 Proven Protease Inhibitor. *Cell* S0092-8674(20)30229–4.
- 311 16. Zhou H, Chen X, Hu T, Li J, Song H, Liu Y, Wang P, Liu D, Yang J, Holmes EC, Hughes  
312 AC, Bi Y, Shi W. 2020. A novel bat coronavirus reveals natural insertions at the S1/S2  
313 cleavage site of the Spike protein and a possible recombinant origin of HCoV-19. *bioRxiv*  
314 2020.03.02.974139.
- 315 17. Lu J, du Plessis L, Liu Z, Hill V, Kang M, Lin H, Sun J, François S, Kraemer MUG, Faria  
316 NR, McCrone JT, Peng J, Xiong Q, Yuan R, Zeng L, Zhou P, Liang C, Yi L, Liu J, Xiao J,  
317 Hu J, Liu T, Ma W, Li W, Su J, Zheng H, Peng B, Fang S, Su W, Li K, Sun R, Bai R, Tang  
318 X, Liang M, Quick J, Song T, Rambaut A, Loman N, Raghvani J, Pybus OG, Ke C. 2020.  
319 Genomic Epidemiology of SARS-CoV-2 in Guangdong Province, China. *Cell*  
320 S0092867420304864.
- 321 18. Quick J, Grubaugh ND, Pullan ST, Claro IM, Smith AD, Gangavarapu K, Oliveira G,  
322 Robles-Sikisaka R, Rogers TF, Beutler NA, Burton DR, Lewis-Ximenez LL, de Jesus JG,

- 323        Giovanetti M, Hill SC, Black A, Bedford T, Carroll MW, Nunes M, Jr LCA, Sabino EC,  
324        Baylis SA, Faria NR, Loose M, Simpson JT, Pybus OG, Andersen KG, Loman NJ. 2017.  
325        Multiplex PCR method for MinION and Illumina sequencing of Zika and other virus  
326        genomes directly from clinical samples. *Nature Protocols* 12:1261–1276.
- 327    19.    Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*  
328        34:3094–3100.
- 329    20.    Grubaugh ND, Gangavarapu K, Quick J, Matteson NL, De Jesus JG, Main BJ, Tan AL,  
330        Paul LM, Brackney DE, Grewal S, Gurfield N, Van Rompay KKA, Isern S, Michael SF,  
331        Coffey LL, Loman NJ, Andersen KG. 2019. An amplicon-based sequencing framework for  
332        accurately measuring intrahost virus diversity using PrimalSeq and iVar. *Genome Biology*  
333        20:8.
- 334    21.    Guindon S, Dufayard J-F, Lefort V, Anisimova M, Hordijk W, Gascuel O. 2010. New  
335        algorithms and methods to estimate maximum-likelihood phylogenies: assessing the  
336        performance of PhyML 3.0. *Syst Biol* 59:307–321.
- 337    22.    Hasegawa M, Kishino H, Yano T. 1985. Dating of the human-ape splitting by a molecular  
338        clock of mitochondrial DNA. *J Mol Evol* 22:160–174.
- 339    23.    Yang Z. 1994. Maximum likelihood phylogenetic estimation from DNA sequences with  
340        variable rates over sites: approximate methods. *J Mol Evol* 39:306–314.
- 341    24.    Wang Y, Song F, Zhu J, Zhang S, Yang Y, Chen T, Tang B, Dong L, Ding N, Zhang Q, Bai Z,  
342        Dong X, Chen H, Sun M, Zhai S, Sun Y, Yu L, Lan L, Xiao J, Fang X, Lei H, Zhang Z,

343 Zhao W. 2017. GSA: Genome Sequence Archive. *Genomics Proteomics Bioinformatics*

344 15:14–18.

345 25. National Genomics Data Center Members and Partners. 2020. Database Resources of the

346 National Genomics Data Center in 2020. *Nucleic Acids Res* 48:D24–D33.

347

348

349 **Figure legends**

350 **Figure 1. Deletion variants identified in SARS-CoV-2 cell strains.** (A) High-throughput sequencing  
351 of the cell-isolated strain (014) from the first SARS-CoV-2 patient (EPI 403934) in Guangdong, China.  
352 Representative reads mapping to the SARS-CoV-2 genome (MN908947.3 used as reference genome)  
353 showed two deletion variants. Redundant proteolytic cleavage sites including furin cleavage site  
354 (PRRARS|V) and cathepsin L site (QSIIAY|T) are marked with red arrows (B) Sanger sequencing of  
355 the 014 cell strains. Heterozygous peaks are highlighted with a red box and sites with distinct three  
356 peaks are marked with \* (C) Results of high-throughput sequencing, showing the ratio of deletion  
357 variants in original clinical sample SF014 (P0) and in cell strains, after 7 rounds of cell passage (P1-7).  
358 The size of square was proportion to the number of reads having these deletions. (D) Phylogenetic tree  
359 of genome sequences of all 24 SARS-CoV-2 cell strains (see Table 1). The size of the circles is  
360 proportional to the percentage of Var1 (QTQTN deletion at 23585–23599) in total reads, except for  
361 strains 619, 4279 and 014 in which Var2 deletions were detected. The maximum likelihood tree was  
362 rooted with the reference genome MN908947.3.

363

364 **Figure 2. The replication kinetics of the deletion variants in Vero-E6 and Vero cells.** Vero-E6 and  
365 Vero cells were infected with the isolated strains 014\_Var1, 014\_Var2, and 029/E6 (Table 1) at  
366 multiplicity of infection (MOI) 0.5. Viral RNA was quantified by real-time PCR using GAPDH as  
367 endogenous control. At the each time point, the relative fold-change in total intracellular viral RNA  
368 was measured by comparison with the viral RNA level at 1-hour post inoculation. Data are the mean  $\pm$   
369 SD of three independent experiments. Asterisk indicate the significant difference ( $p < 0.05$ ).

370

371 **Figure3. Observed deletions near the S1/S2 boundary map to a unresolved region in the cryoEM**  
372 **structure of SARS-CoV-2 S.** Cartoon representation of the SARS-CoV-2 S protein ectodomain, as  
373 resolved by Walls and colleagues(4) (PDB: 6VXX). The S1 and S2 subunits of the different protomers  
374 are indicated (white and grey, respectively). The unresolved loop that contains part of deletion Var1  
375 (<sup>675</sup>QTQTN<sup>679</sup>) and all of deletion Var2 (<sup>679</sup>NSPRRAR<sup>685</sup>) is indicated within each protomer of the  
376 trimeric assembly through signposting flanking residues T<sup>676</sup> and S<sup>689</sup> as spheres in deep teal. Similarly,  
377 the first residue of Var1 (Q<sup>675</sup>), which is resolved in the structure, is indicated as an orange surface  
378 within each of the S protomers. N-linked glycans are shown as blue spheres and the Asn side chains to  
379 which the glycans are linked are presented as sticks. Inset: A zoomed-in side view representation of this  
380 local arrangement is shown. T<sup>676</sup> and S<sup>689</sup>, which flank the unresolved loop, and Var1 residue Q<sup>675</sup> are  
381 numbered and indicated under transparent spheres as deep teal and orange sticks, respectively. A  
382 dashed line indicating the approximate position of the connecting unresolved loop is shown. N-linked  
383 glycans are presented as in the original image with their residue numbers marked.  
384

385 **Table1. Sample information and accession numbers for all sequences**

Patient identifier	Sample isolated from	Passage	Sample name	Sequencing method	Accession number
	BALF	Original	014	Metagenomic	SAMC151281
	Vero-E6	3	014/MiSeq	PCR+MiSeq	SAMC150996
Case1	Vero-E6	3	014/cDNA	Nanopore direct cDNA	SAMC150997
	Vero-E6	Plaque	014_Var1	PCR+Nanopore	SAMC192628
	Vero-E6	Plaque	014_Var2	PCR+Nanopore	SAMC192629
Case2	Vero-E6	2	025/E6	PCR+Nanopore	SAMC150991
	Vero	2	028/Vero	PCR+Nanopore	SAMC150988
Case3	Vero-E6	2	028/E6	PCR+Nanopore	SAMC150992
Case4	Vero-E6	2	029/E6	PCR+Nanopore	SAMC150975
	Vero-E6	2	107/E6	PCR+Nanopore	SAMC150977
Case5	Vero	2	107/Vero	PCR+Nanopore	SAMC150989
	Vero-E6	2	108/E6	PCR+Nanopore	SAMC150993
Case6	Vero	2	108/Vero	PCR+Nanopore	SAMC150995
	Vero-E6	2	112/E6	PCR+Nanopore	SAMC150976
Case7	Vero	2	112/Vero	PCR+Nanopore	SAMC150994
	Vero-E6	2	115/E6	PCR+Nanopore	SAMC150978
Case8	Vero	2	115/Vero	PCR+Nanopore	SAMC150990

Case9	Vero-E6	2	252/E6	PCR+Nanopore	SAMC150980
Case10	Vero-E6	2	262/E6	PCR+Nanopore	SAMC150981
Case11	Vero-E6	2	263/E6	PCR+Nanopore	SAMC150983
Case12	Vero-E6	2	265/E6	PCR+Nanopore	SAMC150982
Case13	Vero-E6	2	272/E6	PCR+Nanopore	SAMC150984
Case14	Vero-E6	3	619/E6	PCR+Nanopore	SAMC153235
Case15	Vero-E6	2	1676/E6	PCR+Nanopore	SAMC150979
Case16	Vero-E6	3	4276/E6	PCR+Nanopore	SAMC153234
Case17	Vero-E6	2	F2/E6	PCR+Nanopore	SAMC150985
Case18	Vero-E6	2	F4/E6	PCR+Nanopore	SAMC150986
Case19	Vero-E6	2	F5/E6	PCR+Nanopore	SAMC150987
Case20	nasopharyngeal	Original	20SF5645	PCR+Nanopore	SAMC150972
Case21	nasopharyngeal	Original	ST-N3-D	PCR+Nanopore	SAMC150973
Case22	nasopharyngeal	Original	SZ-N16-D	PCR+Nanopore	SAMC150974

386

387

388 **Table 2: QTQTN deletion variant (23585–23599, Var1) identified in clinical samples**

Samples	Days post illness onset	REF_depth	ALT_depth	Del Variant Ratio
20SF5645	Asymptomatic	104	25	19.4%
ST-N3-D*	16	82	8	8.8%
SZ-N16-D*	30	256	125	32.8%

389 \* Cases detected with the recurrence of SARS-CoV-2 after discharge





