# Potential transmission chains of variant B.1.1.7 and co-mutations of SARS-CoV-2

Jingsong Zhang[1#], Yang Zhang[2#], Junyan Kang[3,4], Shuiye Chen[2], Yongqun He[5], Benhao Han[1], Mofang Liu[3,4], Lina Lu[1], Li Li[6], Zhigang Yi[2,7]* and Luonan Chen[1,8,9,10]*

[1]State Key Laboratory of Cell Biology, Shanghai Institute of Biochemistry and Cell Biology, Center for Excellence in Molecular Cell Science, Chinese Academy of Sciences, Shanghai 200031, China.

[2]Key Laboratory of Medical Molecular Virology (MOE/NHC/CAMS), School of Basic Medical Sciences, Shanghai Medical College, Fudan University, Shanghai 200032, China.

[3]State Key Laboratory of Molecular Biology, Shanghai Key Laboratory of Molecular Andrology, Shanghai Institute of Biochemistry and Cell Biology, Center for Excellence in Molecular Cell Science, Chinese Academy of Sciences, Shanghai 200031, China.

[4]University of Chinese Academy of Sciences, Shanghai 200031, China.

[5]Department of Computational Medicine and Bioinformatics, University of Michigan Medical School, Ann Arbor, MI 48109, USA.

[6]Department of Genetics, Harvard Medical School, Boston, MA 02115, USA.

[7]Shanghai Public Health Clinical Center, Fudan University, Shanghai 201508, China.

[8]School of Life Science and Technology, ShanghaiTech University, Shanghai 201210, China.

[9] Key Laboratory of Systems Health Science of Zhejiang Province, Hangzhou Institute for Advanced Study, University of Chinese Academy of Sciences, Hangzhou 310024, China.

[10] Pazhou Lab, Guangzhou 510330, China.

*Correspondence: Luonan Chen (lnchen@sibs.ac.cn) or Zhigang Yi (zgyi@fudan.edu.cn)

#These authors contributed equally: Jingsong Zhang, Yang Zhang

1

# **Abstract**

The presence of SARS-CoV-2 mutants, including the emerging variant B.1.1.7, has raised great concerns in terms of pathogenesis, transmission, and immune escape. Characterizing SARS-CoV-2 mutations, evolution, and effects on infectivity and pathogenicity is crucial to the design of antibody therapies and surveillance strategies. Here we analyzed 454,443 SARS-CoV-2 spike genes/proteins and 14,427 whole-genome sequences. We demonstrated that the early variant B.1.1.7 may not have evolved spontaneously in the United Kingdom or within human populations. Our extensive analyses suggested that Canidae, Mustelidae or Felidae, especially the Canidae family (for example, dog) could be a possible host of the direct progenitor of variant B.1.1.7. An alternative hypothesis is that the variant was simply yet to be sampled. Notably, the SARS-CoV-2 whole genome represents a large number of potential co-mutations with very strong statistical significances (p value<E–44). In addition, we used an experimental SARS-CoV-2 reporter replicon system to introduce the dominant co-mutations NSP12_c14408t, 5'UTR_c241t, and NSP3_c3037t into the viral genome, and to monitor the effect of the mutations on viral replication. Our experimental results demonstrated that the co-mutations significantly attenuated the viral replication. The study provides valuable clues for discovering the transmission chains of variant B.1.1.7 and understanding the evolutionary process of SARS-CoV-2.

**Key words:** SARS-CoV-2, variant B.1.1.7, transmission chains, co-mutations, viral replication.

## Introduction

Since the outbreak in December 2019, COVID-19 has been pandemic in over 200 countries. Cases of infection and mortalities have been surging and are an ongoing threat to public health[1,2]. COVID-19 is caused by infection with the novel coronavirus SARS-CoV-2[3-5]. Although as a coronavirus, SARS-CoV-2 has genetic proofreading mechanisms[6-8], the persistent natural selection pressure in the population drives the virus to gradually accumulate favorable mutations[6,9,10]. Much attention has been paid to the mutations and evolution of SARS-CoV-2[11-15], since mutations are related to the infectivity and pathogenicity of viruses[16-21]. Beneficial mutants of the virus can better evolve and adapt to the host[9], either strengthening or weakening the infectivity and pathogenicity. In addition, certain variants may generate drug resistance and reduce the efficacy of vaccines and therapeutics[22-26]. In short, studying mutations and evolution in detail is vital to understand the transformations of viral properties and to control the pandemic.

A new variant of SARS-CoV-2 named VOC-202012/01 (Variant of Concern 202012/01) or lineage B.1.1.7 was first detected in the United Kingdom last December[27]. It appears to be substantially more transmissible than other variants[28]. The variant has been growing exponentially in the United Kingdom and rapidly spreading to other countries[29,30]. However, it is not yet clear if it evolved spontaneously in the United Kingdom or was imported from other countries. Studying how the variant B.1.1.7 mutates can enable researchers to track its spread over time and to understand the evolution of SARS-CoV-2.

In this study, large-scale SARS-CoV-2 sequences, consisting of more than 454,000 spike

genes/proteins and 14,000 whole-genome sequences were analyzed. Our extensive sequence analysis showed that many mutations always co-occur not only in the spike protein of B.1.1.7, but in the whole genome of SARS-CoV-2. The mutation trajectories of the spike protein indicate that the early variant B.1.1.7 did not evolve spontaneously in the United Kingdom or even within human populations. We also investigated possible SARS-CoV-2 transmission chains of the variant B.1.1.7 based on the mutation analysis of large-scale spike proteins and the cluster analysis of spike genes. Over the whole genome, the top 25 high-frequency mutations of SARS-CoV-2 converged into several potential co-mutation patterns, each of which showed a strong correlation with a very strong statistical significance (p value<E–44). The potential co-mutations depicted the evolutionary trajectory of SARS-CoV-2 virus in the population, shaping variable replication of SARS-CoV-2. In addition, we further explored the effect of the dominant (co-)mutations 5'UTR_c241t, NSP3_c3037t, and NSP12_c14408t on viral replication using a SARS-CoV-2 replicon based on a four plasmid *in-vitro* ligation system. The results suggest that such mutations significantly attenuate the replication of SARS-CoV-2.

## Results

## Evolutionary trajectories of variant B.1.1.7

The variant B.1.1.7 was generally defined by multiple amino acid changes including 3 deletions (69-70del and 145del) and 7 mutations (N501Y, A570D, D614G, P681H, T716I, S982A, and D1118H) in the spike protein[31]. The number of non-adjacent co-occurrent changes indicates that they resulted from accumulated mutations. We therefore explored the

4

evolutionary trajectories of B.1.1.7 by tracing the incremental mutations (Fig. 1a). All routes along the directions of the arrows are possible evolutionary trajectories of lineage B.1.1.7. Among all the mutation routes, the green one was the most probable mutation trajectory based on the number of variant strains. However, it was unlikely that the earliest variant B.1.1.7 (GISAID: EPI_ISL_601443, 2020-09-20, England) with 9 mutations evolved from the existing variants with 3–8 mutations, because the former arose much earlier than the latter. More than 454,000 SARS-CoV-2 strains have been collected and extensive sequenced from infected humans without finding intermediate variants with 3–9 mutations. It is therefore unlikely that the intermediate variants with 3–8 mutations have infected humans. Thus, the early variant B.1.1.7 might not have arisen spontaneously in the UK or within human populations. An alternative hypothesis is that spillover likely occurred from susceptible animals.

The co-appearance rates (see Materials and Methods) of all nine mutations are shown in Fig. 1b. We found that at least five mutations (145del, A570D, T716I, S982A, and D1118H) of variant B.1.1.7 significantly co-occurred (rate>95%), which indicates a potential co-mutation pattern in the spike protein, causing us to wonder what selection pressure drove such co-occurrences of mutations and rapid evolution in the population of SARS-CoV-2. Note that coronaviruses generally tend to exhibit rapid evolution when they jump to a different species[32]. We therefore analyzed the key spike genes and proteins of existing SARS-CoV-2 strains collected from animals to find a possible direct progenitor of variant B.1.1.7. The variant with mutations "56" (labeled by "*" in Fig. 1a, termed star variant) had the minimum phylogenetic distance with EPI_ISL_699508, which was collected from a dog

109  on 2020-07-28 (Fig. 2) using MEGA[33,34] (see Materials and Methods). The strains collected

110  from tigers, minks, and cats were also close to the star variant. Our extensive analyses

111  including mutations, phylogeny (Fig. 2), collection date/location and the number of

112  sequences (Tables S1-S3) suggested that Canidae, Mustelidae or Felidae, especially the

113  Canidae family (for example, dog) could be a possible host of the direct progenitor of variant

114  B.1.1.7. The possible transmission chains of variant B.1.1.7 are shown in Fig. 1c. This star

115  variant strains in humans could not have evolved into the early variant B.1.1.7, but they

116  might have infected high-density yet susceptible animals (such as dogs) and adapted to

117  these species through rapid mutation. Such progenitor variants comprised most or all of the

118  mutations of the early variant B.1.1.7 within the Canidae family populations, and they may

119  have spilled back to humans after the rapid mutation period.

## High-frequency mutations converge into potential co-mutations

121  Based on sequence alignment and mutation analysis, we found that 7,441 nucleotide

122  alterations in the viral 29903-letter RNA code occurred at least once in the samples from

123  COVID-19 patients. These mutations were dispersed in the 14,427 SARS-CoV-2 strains

124  collected from all around the world. As shown in the heatmap of the top 1% high-frequency

125  mutations (Table S4), some sites show very similar mutation rates on most days in samples

126  isolated globally (Fig. S1), including 8,898 and 815 samples isolated from the U.S. (Fig. S2)

127  and Australia (Fig. S3). Therefore, these mutations shown in Fig. S4a were selected and

128  clustered into co-occurrences, which we called potential co-mutation patterns. From the

129  landscape of the mutation rates (Fig. S4a), 25 nucleotide sites were clearly clustered into

130 several potential co-mutation patterns. Among these patterns, there was one consisting of

131 the top 4 high-frequency mutations (i.e., 5'UTR_c241t, NSP3_c3037t, NSP12_c14408t, and

132 S_a23403g), which converged into a dominant potential co-mutation pattern. Such

133 co-occurrence lineage has been found in almost all sequenced samples of SARS-CoV-2.

134 Within this co-occurrence pattern, mutation S_ a23403g resulted in the amino acid change

135 (D614G) that apparently enhances viral infectivity[6,35], albeit debate exists[16]. Notably, there

136 were three successive sites at the 28881$^{st}$ to 28883$^{rd}$ positions of the virus (N_g28881a,

137 N_g28882a, and N_g28883c) that strictly co-occurred. Comparing Fig. S4a-c and Table S4,

138 we found that the top 14 high-frequency mutations formed five common co-occurrence

139 patterns.

140 To assess the above co-occurrence patterns, we analyzed the correlations and statistical

141 significance levels of the high-frequency co-occurrence mutations. The heatmap of the

142 paired Pearson-correlation-coefficients (Fig. 3a) shows that the top 25 high-frequency

143 mutations clearly cluster into several potential co-mutation groups/patterns with very strong

144 correlation (≥0.8). By regression analyses, the above co-occurrence patterns have statistical

145 significance levels with p values less than 10$^{-44}$ (Fig. 3b). The detailed mutation transitions

146 (Fig. 3c–k, Figs. S5–7) provide further evidence that the above mutations form co-mutation

147 patterns.

## Dominant mutations attenuate viral replication

149 We further explored the effect of the dominant mutations 5'UTR_c241t, NSP3_c3037t,

150 and NSP12_c14408t on viral replication using a SARS-CoV-2 replicon based on a

151 four-plasmid *in-vitro* ligation system. This replicon is devoid of the viral structural proteins

152 while undergoing viral replication, and the viral replication is sensitive to the antiviral agent

153 remdesivir[36]. The 5'UTR_c241t mutation resides in a highly conserved region in the 5'UTR

154 (Fig. 4a). The NSP3_c3037t mutation is synonymous. The NSP12_c14408t mutation is

155 nonsynonymous with an amino acid change of a conserved amino acid P323 in the viral

156 RNA-dependent RNA polymerase (Fig. 4b). We introduced the NSP12_c14408t mutation or

157 the NSP12_c14408t mutation with the other two mutations 5'UTR_c241t and NSP3_c3037t

158 into the replicon plasmids. The fragments were released from the plasmids by BsaI digestion,

159 and then assembled by *in-vitro* ligation with T4 ligase (Fig. 4c). Replicon RNA transcribed

160 from the ligation products was co-transfected with N mRNA into Huh7 cells. RNA replication

161 was monitored by measuring the secreted *Gaussia* luciferase activity in the supernatants.

162 Enzymatic dead mutants (759-SAA-761) of the RNA-dependent RNA polymerase NSP12

163 were introduced, and the mutated replicon served as a non-replication control. As shown in

164 Fig. 4d, transfection of WT replicon RNA resulted in an obvious increase of luciferase activity,

165 and SAA RNA did not replicate as expected. Introduction of NSP12_c14408t mutation

166 resulted in a significant reduction of viral replication. The combination of NSP12_c14408t

167 mutation with the other two mutations further significantly but only marginally reduced viral

168 replication. These results demonstrate that the P323L mutation in the viral RNA-dependent

169 RNA polymerase reduces viral replication, and the synonymous mutations may further

170 attenuate viral replication.

# Discussion

A well-resolved phylogeny of variant B.1.1.7 spike genes provides an opportunity to understand the evolutionary process and transmission chains of variant B.1.1.7. Our incremental mutation and phylogenetic analyses on large-scale SARS-CoV-2 spike proteins/genes revealed that the early variant B.1.1.7 might not have evolved spontaneously in the United Kingdom or within human populations. In this case the spillover likely occurred from susceptible animals. Current evidence[37-39] indicates that SARS-CoV-2 can effectively infect both domestic animals (for example, dog, cat, pig and bovine) and wild animals (for example, mink, rabbit and fox) by binding their angiotensin converting enzyme 2 (ACE2). Our further analyses including mutations, phylogeny, collection date/location and the number of sequences suggested that the earliest variant B.1.1.7 possibly originated from Canidae, Mustelidae or Felidae, especially the Canidae family (for example, dog). The cases[40] that the variant B.1.1.7 can easily infect dogs and cats indicated that both are susceptible to B.1.1.7. Still, due to the limited information available to date, an alternative hypothesis is that the direct progenitor of variant B.1.1.7 is yet to be sampled. In addition to variant B.1.1.7, as a future topic we will work on the analysis of other lineages such as P.1, B.1.351, B.1.427, and B.1.42, when sufficient numbers of their sequences are available.

By tracing the mutation trajectories, we found that at least five mutations of the spike proteins always co-occurred, and a large number of potential co-mutations appeared in the top 1% high-frequency mutations of SARS-CoV-2 whole genome. It has been documented that the mutation S_ a23403g results in the amino acid change of the spike protein D614G

192 and enhances viral infectivity[19,41-44]. Here, by using a SARS-CoV-2 reporter replicon system,

193 we demonstrated that the one of the dominant co-mutations NSP12_c14408t significantly

194 reduced viral replication and combination of NSP12_c14408t mutation with the other two

195 synonymous mutations 5'UTR_c241t and NSP3_c3037t although significantly but only

196 marginally reduced viral replication further. As the 5'UTR play an important role in regulating

197 viral replication, the synonymous mutations 5'UTR_c241t may attenuate viral replication by

198 change RNA secondary structure[45]. These findings imply that SARS-CoV-2 undergoes an

199 evolution toward enhancing viral infectivity while attenuating viral replication. SARS-CoV-2

200 has exhibited significant mutations and co-mutations. We evaluated the replication of a

201 co-mutation pattern including three dominant mutations. If other mutations act similarly on

202 the viral replication needs to be verified. These results can be further explored for efficient

203 vaccine design in our future work. In summary, this study provides insights into the

204 transmission chains of variant B.1.1.7 and the effect of viral dominant mutations on viral

205 evolution.

## Materials and Methods

206

### Data selection and pre-processing

207

208 The 454,443 spike gene/protein sequences of SARS-CoV-2 were obtained at

209 https://www.gisaid.org/. The NCBI website at https://www.ncbi.nlm.nih.gov/sars-cov-2/ has

210 released more than 1.7 thousand sequences of SARS-CoV-2 viruses before July 31, 2020.

211 We selected 14,427 sequences that satisfied two criteria: (1) having specific collection dates;

212 (2) sequence-lengths being no less than 29,305 nt (29903*0.98). It is inevitable that some

213 sites of sequences are equivocal owing to the limitation of sequencing depth. For instance,

214 many sites were labeled as letter N in genome sequences. The noise of indeterminate

215 nucleic-acids was taken into consideration in our experiments so as to boost accuracy. The

216 co-mutation rate of multi-site co-mutations was calculated by $co-\text{mutation rate} =$

217 $\frac{\text{number of sequences containing co}-\text{mutaions}}{\text{number of all sequences}}$. Moreover, the co-appearance rate of a mutation in

218 B.1.1.7 variant was defined by $co-\text{appearance rate} = \frac{\text{number of B.1.1.7 sequences}}{\text{number of sequences containing a mutation}}$.

**Possible animal host analyses**

220 In addition to the phylogenetic analysis, we further explored the possible animal hosts of

221 the direct progenitor of variant B.1.1.7 by mutations, collection time/space of strains, the

222 number of sequences and the edit distance[46,47] of mutations (Table S1-2). Due the late

223 lockdown policies of some governmental agencies, the spread of SARS-CoV-2 has not been

224 prevented well in Europe, America, and Australia. We could ignore the impact of policies for

225 studying the origin of variant B.1.1.7. We quantified the multiple impact factors of viral

226 transmission as shown in Table S3 based on the criterion that the smaller the value, the

227 more similar. The results still supported that the Canidae family is a possible host of the

228 direct progenitor of variant B.1.1.7.

**MEGA version and parameter settings**

230 Version: MEGA-X

231 Statistical Method: Maximum Likelihood

232 Test of Phylogeny: None

233 Model/Method: Jones-Taylor-Thornton (JTT) model

234    Rates among Sites: Uniform Rates

235    Gaps/Missing Data Treatment: Use all sites

236    ML Heuristic Method: Nearest-Neighbor-Interchange (NNT)

237    Initial Tree for ML: Make initial tree automatically (Default - NJ/BioNJ)

238    Branch Swap Filter: None

239    Number of Threads: 7

## Statistical analysis

241    The Pearson-correlation-coefficient (PCC) is a classic statistic that measures linear

242    correlation between two variables. Its value ranges from -1.0 to 1.0. Normally, the two

243    variables meet a strong correlation or a very strong correlation when the absolutes value of

244    PCC is between 0.6 and 0.8 or between 0.8 and 1.0. Linear regression is a linear approach

245    to model the relationship between a scalar response and one or more variables. We used

246    PCC and significance level (p value) of regression analysis to evaluate the relationships of

247    the co-occurrence mutations in large-scale SARS-CoV-2 examples.

## Plasmids

249    Four plasmids encompassing the viral genome (pLC-nCoV-A-BsaI, pLC-nCoV-B-BsaI,

250    pLC-nCoV-C-BsaI, and pnCoV-D-sGluc-BsaI) were described previously[36]. The

251    5'UTR_c-241-t and NSP3_c-3037-t mutations were introduced into the pLC-nCoV-A-BsaI by

252    fusion PCR. The NSP12_c-14408-t mutation was introduced into the pLC-nCoV-B-BsaI by

253    fusion PCR.

## Cell lines

12

255    The human hepatoma cells Huh 7 were purchased from the Cell Bank of the Chinese

256    Academy of Sciences (www.cellbank.org.cn) and routinely maintained in Dulbecco's

257    modified medium supplemented with 10% FBS (Gibco) and 25 mM HEPES (Gibco).

### *In-vitro* ligation

259    BsaI digested fragments were gel purified using Gel Extraction Kit (OMEGA) and ligated

260    with T4 ligase (New England Biolabs) at room temperature for 1 h. The ligation products

261    were phenol/chloroform extracted, precipitated by absolute ethanol, and resuspended in

262    nuclease-free water, quantified by determining the A260 absorbance.

### *In-vitro* transcription

264    Purified *in-vitro* ligated product was used as template for the *in-vitro* transcription by

265    mMESSAGE mMACHINE T7 Transcription Kit (Ambion) according to the manufacturer's

266    protocol. For N mRNA production, we amplified the N coding region by PCR (sense: *GGC*

267    *ACA CCC CTT TGG CTC T*; antisense: *TTT TTT TTT TTT TTT TTT TTT TTT TTT TTT TTT*

268    *TTT TTT TCT AGG CCT GAG TTG AGT CAG CAC*) with phCMV-N as template. Then the

269    purified PCR product was used as a template for *in-vitro* transcription by mMESSAGE

270    mMACHINE T7 Transcription Kit as described above. RNA was purified by RNeasy mini

271    Elute (Qiagen), eluted in nuclease-free water, and quantified by UV absorbance (260 nm).

### Transfection

273    Cells were seeded onto 48-well plates at a density of $7.5 \times 10^4$ per well and then

13

274 transfected with 0.3 μg *in-vitro* transcribed RNA using a TransIT-mRNA transfection kit

275 (Mirus) according to the manufacturer's protocol.

276 **Luciferase activity**

277    Supernatants were taken from cell medium and mixed with equal volumes of 2×lysis

278 buffer (Promega). Luciferase activity was measured with Renilla luciferase substrate

279 (Promega) according to the manufacturer's protocol.

# Acknowledgments

# Author contributions

292 L.N.C. and J.S.Z. designed the study. Z.G.Y. and J.S.Z. designed the experiments. J.S.Z.

293 analyzed data. Y. Z. performed the experiments of viral replication. J.S.Z., Z.G.Y., and J.Y.K.

294 designed the figures. S.C. repeated and checked the experiments of viral replication. H.B.H.

295 checked the computational analyses. J.S.Z. and Z.G.Y. wrote the manuscript. Y.Q.H, M.F.L.,

296 L.N.L, and L.L. polished the manuscript. All authors participated in result interpretation and

297 discussion.

## Data availability

299 The raw sequence data reported in this paper have been deposited in the GISAID and NCBI

300 websites    at    https://www.gisaid.org/    and    https://www.ncbi.nlm.nih.gov/sars-cov-2/,

301 respectively. Code is available from the corresponding author on reasonable request.

## Conflict of interest

303 The authors declare that they have no conflict of interest.

## References

305 1    Wrapp, D. *et al.* Cryo-EM structure of the 2019-nCoV spike in the prefusion conformation. *Science* **367**,
306      1260-1263 (2020).
307 2    Chinazzi, M. *et al.* The effect of travel restrictions on the spread of the 2019 novel coronavirus (COVID-19)
308      outbreak. *Science*, 395-400 (2020).
309 3    Zhu, N. *et al.* A Novel Coronavirus from Patients with Pneumonia in China, 2019. *The New England journal of*
310      *medicine* **382**, 727-733, doi:10.1056/NEJMoa2001017 (2020).
311 4    Zhou, P. *et al.* A pneumonia outbreak associated with a new coronavirus of probable bat origin. *Nature* **579**,
312      270-273, doi:10.1038/s41586-020-2012-7 (2020).
313 5    Chen, L. *et al.* RNA based mNGS approach identifies a novel human coronavirus from two individual pneumonia
314      cases    in    2019    Wuhan    outbreak.    *Emerging    microbes    &    infections* **9**,    313-319,
315      doi:10.1080/22221751.2020.1725399 (2020).
316 6    Korber, B. *et al.* Tracking changes in SARS-CoV-2 Spike: evidence that D614G increases infectivity of the
317      COVID-19 virus. *Cell*, doi:doi.org/10.1016/j.cell.2020.06.043 (2020).
318 7    Smith, E. C., Blanc, H., Vignuzzi, M. & Denison, M. R. Coronaviruses Lacking Exoribonuclease Activity Are
319      Susceptible to Lethal Mutagenesis: Evidence for Proofreading and Potential Therapeutics. *Plos Pathog* **9** (2013).

15

8    Sevajol, M., Subissi, L., Decroly, E., Canard, B. & Imbert, I. Insights into RNA synthesis, capping, and proofreading mechanisms of SARS-coronavirus. *Virus Res* **194**, 90-99 (2014).

9    Vignuzzi, M., Stone, J. K., Arnold, J. J., Cameron, C. E. & Andino, R. J. N. Quasispecies diversity determines pathogenesis through cooperative interactions in a viral population.   **439**, 344-348 (2006).

10   Garvin, M. R. *et al.* Potentially adaptive SARS-CoV-2 mutations discovered with novel spatiotemporal and explainable AI models. *Genome Biology* **21**, 304, doi:10.1186/s13059-020-02191-0 (2020).

11   Tang, X. *et al.* On the origin and continuing evolution of SARS-CoV-2. *National Science Review* **7**, 1012-1023, doi:10.1093/nsr/nwaa036 %J National Science Review (2020).

12   Acter, T. *et al.* Evolution of severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) as coronavirus disease 2019 (COVID-19) pandemic: A global health emergency. *Science of The Total Environment* **730**, 138996, doi:https://doi.org/10.1016/j.scitotenv.2020.138996 (2020).

13   Plante, J. A. *et al.* Spike mutation D614G alters SARS-CoV-2 fitness. *Nature*, doi:10.1038/s41586-020-2895-3 (2020).

14   Zohar, T. *et al.* Compromised Humoral Functional Evolution Tracks with SARS-CoV-2 Mortality. *Cell* **183**, 1508-1519.e1512, doi:10.1016/j.cell.2020.10.052 (2020).

15   Li, T. *et al.* The use of SARS-CoV-2-related coronaviruses from bats and pangolins to polarize mutations in SARS-Cov-2. *Science China Life Sciences* **63**, 1608-1611, doi:10.1007/s11427-020-1764-2 (2020).

16   Grubaugh, N. D., Hanage, W. P. & Rasmussen, A. L. Making sense of mutation: what D614G means for the COVID-19 pandemic remains unclear. *Cell*, doi:doi.org/10.1016/j.cell.2020.06.040 (2020).

17   Liu, Z. *et al.* Identification of Common Deletions in the Spike Protein of Severe Acute Respiratory Syndrome Coronavirus 2. *J Virol* **94** (2020).

18   Li, Q. Q. *et al.* The Impact of Mutations in SARS-CoV-2 Spike on Viral Infectivity and Antigenicity. *Cell* **182**, 1284-+ (2020).

19   Plante, J. A. *et al.* Spike mutation D614G alters SARS-CoV-2 fitness. *Nature* (2020).

20   Blanco, J. D., Hernandez-Alias, X., Cianferoni, D. & Serrano, L. In silico mutagenesis of human ACE2 with S protein and translational efficiency explain SARS-CoV-2 infectivity in different species. *Plos Comput Biol* **16** (2020).

21   Zhang, L. Z. *et al.* SARS-CoV-2 spike-protein D614G mutation increases virion spike density and infectivity. *Nat Commun* **11** (2020).

22   A., B. *et al.* Antibody cocktail to SARS-CoV-2 spike protein prevents rapid mutational escape seen with individual antibodies. *Science*, doi:10.1126/science.abd0831 (2020).

23   Hansen, J. *et al.* Studies in humanized mice and convalescent humans yield a SARS-CoV-2 antibody cocktail. *Science*, doi:10.1126/science.abd0827 (2020).

24   Sheahan, T. P. *et al.* An orally bioavailable broad-spectrum antiviral inhibits SARS-CoV-2 in human airway epithelial cell cultures and multiple coronaviruses in mice. *Sci Transl Med* **12** (2020).

25   Nunes-Santos, C. J., Kuehn, H. S. & Rosenzweig, S. D. N-Glycan Modification in Covid-19 Pathophysiology: In vitro Structural Changes with Limited Functional Effects. *J Clin Immunol* (2020).

26   Lo, M. K. *et al.* Remdesivir targets a structurally analogous region of the Ebola virus and SARS-CoV-2 polymerases. *P Natl Acad Sci USA* **117**, 26946-26954 (2020).

27   Editorial. Evolution goes viral. *Nature Ecology & Evolution* **5**, 143-143, doi:10.1038/s41559-021-01395-2 (2021).

28   Editorial. COVID-19 vaccines: acting on the evidence. *Nature Medicine*, doi:10.1038/s41591-021-01261-5 (2021).

29   Muik, A. *et al.* Neutralization of SARS-CoV-2 lineage B.1.1.7 pseudovirus by BNT162b2 vaccine–elicited human sera. eabg6105, doi:10.1126/science.abg6105 %J Science (2021).

364  30  Rice, B. L. *et al.* Variation in SARS-CoV-2 outbreaks across sub-Saharan Africa. *Nature Medicine*,
365      doi:10.1038/s41591-021-01234-8 (2021).

366  31  Muik, A. *et al.* Neutralization of SARS-CoV-2 lineage B.1.1.7 pseudovirus by BNT162b2 vaccine–elicited human
367      sera. *Science*, eabg6105, doi:10.1126/science.abg6105 %J Science (2021).

368  32  Zhou, P. & Shi, Z.-L. SARS-CoV-2 spillover events. *Science* **371**, 120-122, doi:10.1126/science.abf6097 %J Science
369      (2021).

370  33  Kumar, S., Stecher, G., Li, M., Knyaz, C. & Tamura, K. MEGA X: Molecular Evolutionary Genetics Analysis across
371      Computing Platforms. *Molecular Biology and Evolution* **35**, 1547-1549 (2018).

372  34  Tamura, K. & Nei, M. Estimation of the Number of Nucleotide Substitutions in the Control Region of
373      Mitochondrial-DNA in Humans and Chimpanzees. *Molecular Biology and Evolution* **10**, 512-526 (1993).

374  35  Li, Q. *et al.* The impact of mutations in SARS-CoV-2 spike on viral infectivity and antigenicity. *Cell*,
375      doi:doi.org/10.1016/j.cell.2020.07.012 (2020).

376  36  Zhang, Y., Song, W., Chen, S., Yuan, Z. & Yi, Z. A bacterial artificial chromosome (BAC)-vectored noninfectious
377      replicon of SARS-CoV-2. *bioRxiv*, doi:doi.org/10.1101/2020.09.11.294330 (2020).

378  37  Shi, J. *et al.* Susceptibility of ferrets, cats, dogs, and other domesticated animals to SARS–coronavirus 2. *Science*
379      **368**, 1016-1020, doi:10.1126/science.abb7015 %J Science (2020).

380  38  Wu, L. *et al.* Broad host range of SARS-CoV-2 and the molecular basis for SARS-CoV-2 binding to cat ACE2. *Cell*
381      *Discovery* **6**, 68, doi:10.1038/s41421-020-00210-9 (2020).

382  39  Patterson, E. I. *et al.* Evidence of exposure to SARS-CoV-2 in cats and dogs from households in Italy. *Nat*
383      *Commun* **11**, 6231, doi:10.1038/s41467-020-20097-0 (2020).

384  40  Ferasin, L. *et al.* Myocarditis in naturally infected pets with the British variant of COVID-19. *bioRxiv*,
385      2021.2003.2018.435945, doi:10.1101/2021.03.18.435945 %J bioRxiv (2021).

386  41  Raghav, S. *et al.* Analysis of Indian SARS-CoV-2 Genomes Reveals Prevalence of D614G Mutation in Spike
387      Protein Predicting an Increase in Interaction With TMPRSS2 and Virus Infectivity. *Front Microbiol* **11** (2020).

388  42  Johnson, M. C. *et al.* Optimized Pseudotyping Conditions for the SARS-COV-2 Spike Glycoprotein. *J Virol* **94**
389      (2020).

390  43  Jiang, X. Y. *et al.* Bimodular effects of D614G mutation on the spike glycoprotein of SARS-CoV-2 enhance
391      protein processing, membrane fusion, and viral infectivity. *Signal Transduct Tar* **5** (2020).

392  44  Fernandez, A. Structural Impact of Mutation D614G in SARS-CoV-2 Spike Protein: Enhanced Infectivity and
393      Therapeutic Opportunity. *Acs Medicinal Chemistry Letters* **11**, 1667-1670 (2020).

394  45  Sun, L. *et al.* <em>In vivo</em> structural characterization of the SARS-CoV-2 RNA genome identifies
395      host proteins vulnerable to repurposed drugs. *Cell* **184**, 1865-1883.e1820, doi:10.1016/j.cell.2021.02.008
396      (2021).

397  46  Ristad, E. S., Yianilos, P. N. J. I. T. o. P. A. & Intelligence, M. Learning string-edit distance. *IEEE Transactions on*
398      *Pattern Analysis and Machine Intelligence* **20**, 522-532 (1998).

399  47  Bille, P. J. T. c. s. A survey on tree edit distance and related problems. *Pattern Analysis and applications* **337**,
400      217-239 (2005).

## Figures and Figure legends

401

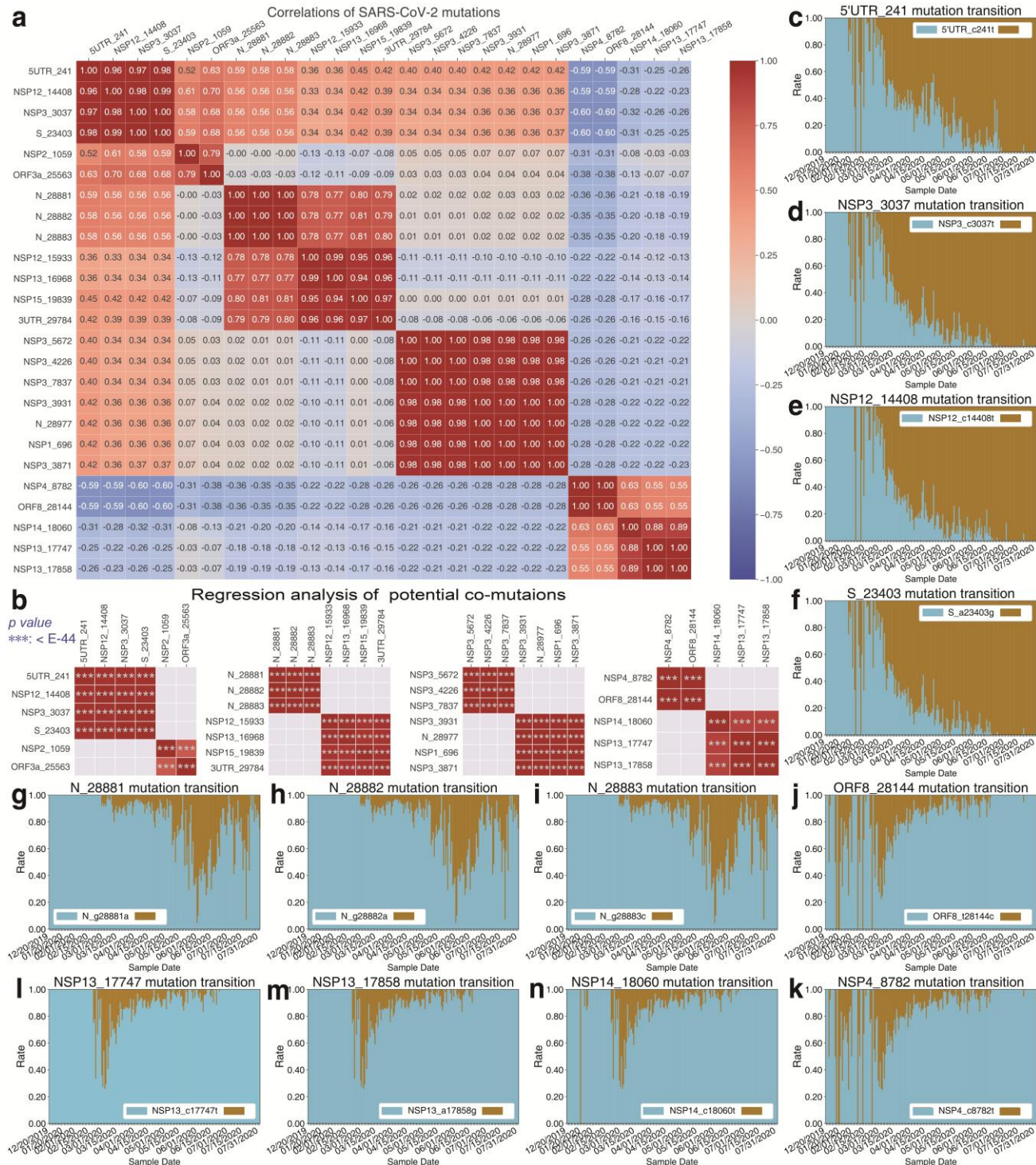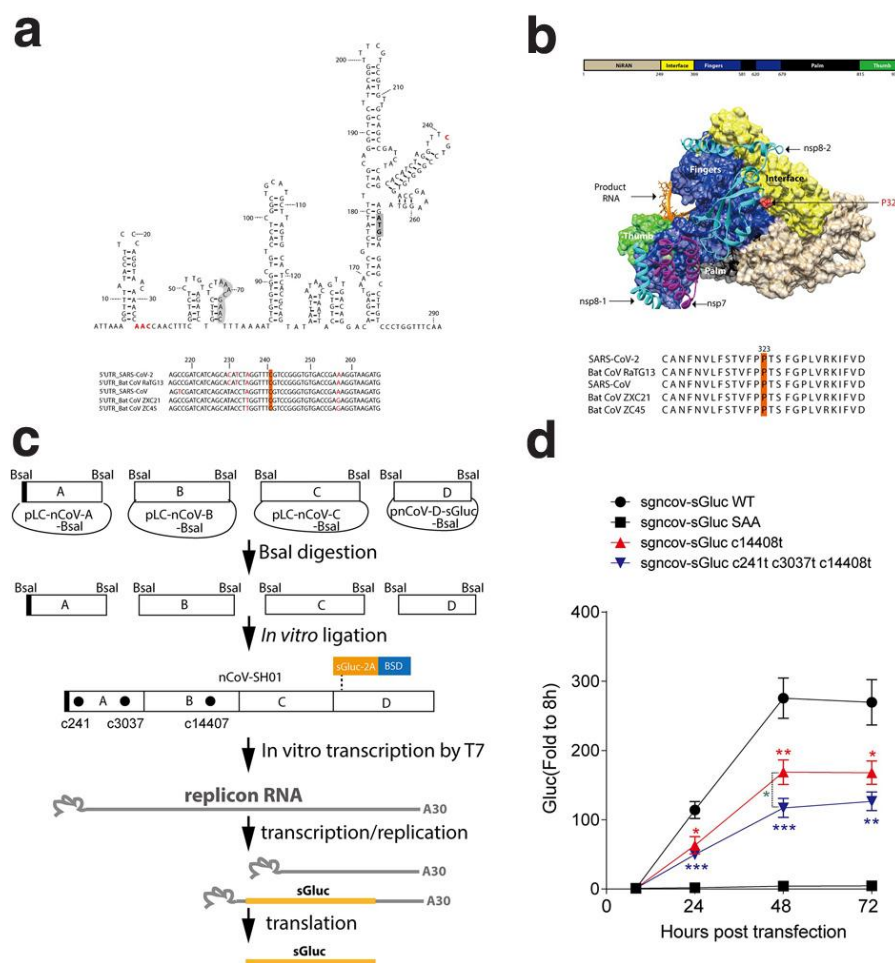### a      Evolutionary trajectories of variant B.1.1.7



402

403 **Fig. 1 Evolutionary trajectories of variant B.1.1.7. a** Incremental mutations of variant B.1.1.7. The

404 digits in the upper-right-corner rectangle with dotted line indicate the labels of mutations. For

405 simplicity, the 69–70 deletions were labeled as "1", and the other mutations "2"-"9" respectively. The

406 bottom nodes (rectangles) represent the variants with one mutation and the top one was the early

407 variant B.1.1.7. Each rectangle with solid line consists of lineage (e.g., B.1.243), number of strains

408 (e.g., N:2382), mutation sites (e.g., M:----56---), the earliest collection date (e.g., 20-03-29, i.e.,

409 2020-03-29), and collection location (e.g., USA). In the labels of the mutation sites, sign "-" indicated

410 the corresponding site did not mutate. All routes along the directions of the arrows are possible

411 evolutionary trajectories of lineage B.1.1.7, where the green one was the most probable mutation

412 trajectory. Large-scale SARS-CoV-2 analysis demonstrates that the early variant B.1.1.7 might not

413 have arisen spontaneously in the UK or within human hosts. **b** Coappearances of variant B.1.1.7

414 mutations. At least five mutations form a potential co-mutation pattern (coappearance rate > 95%). **c**

415 Possible transmission chains of variant B.1.1.7. Canidae, Mustelidae or Felidae, especially the

416 Canidae family (for example, dog) could be a possible host of the direct progenitor of variant B.1.1.7.

417

## Phylogenetic analysis of variant B.1.1.7 by Spike

**418**

**419** **Fig. 2 The Canidae family could be a possible host of the direct progenitor of variant B.1.1.7.**

**420** The digits on the left of the figure indicate the labels of mutations, which correspond with the

**421** mutation labels in Fig. 1a. The strains shown in the center of the figure contain at least one spike

**422** mutation of variant B.1.1.7. And these strain examples cover all existing SARS-CoV-2 viruses that

**423** collected from animal hosts. The strain labeled by orange star corresponds with the star variant in Fig.

**424** 1a. The strain with orange solid-round label was collected from a dog on 2020-07-28. Such two

**425** strains share the same mutations "56" and have the minimum phylogenetic distance by MEGA tool.

**426** Canidae, Mustelidae or Felidae, especially the Canidae family (for example, dog) could be a possible

**427** host of the direct progenitor of variant B.1.1.7 based on existing stains collected before the end of

**428** Jan. 2021.

**Fig. 3 The strong correlations suggest that the top 25 mutations form eight potential co-mutation patterns. a** The correlation heatmap of the top 25 mutations. These mutations could be grouped into several clusters with high Pearson-correlation-coefficient (PCC). **b** Regression analysis of mutations shows that eight clusters all denote the statistical significance level: ***p value < E-44. c to k show the transitions of the high-frequency mutations. The sky-blue represents the rate per day of initial residue in population and the golden the rate per day of substitution/mutant. These mutation transitions provide further evidence that the above mutations potentially form co-mutation patterns.

**Fig. 4 Dominant co-mutation attenuates viral replication. a** Predicted RNA structure of the SARS-CoV-2 5'UTR. RNA structure of the 400-nt 5'UTR was predicted by "RNAstructure" (http://rna.urmc.rochester.edu/RNAstructureWeb). The start codon for nsp1 is grey, the TRS-L is orange, and the mutated nucleotides are red. The bottom panel shows the alignment of the 5'UTR of SARS-CoV-2 with 5'UTRs of related viruses, with c241 highlighted. **b** Structure of SARS-CoV-2 RdRp/RNA complex. The structure of SARS-CoV-2 RdRp/RNA complex (PDB, 6X2G) was visualized by Chimera (UCSF). The P323 mutation is highlighted in red, with the alignment of the amino acid sequences of SARS-CoV-2 and related viruses near the P323 position. **c** Schematic of the *in-vitro* ligation system for SARS-CoV-2 replicon. Four plasmids encompassing the viral genome were digested by BsaI to release the four fragments. After gel purification, the fragments were ligated by T4 ligase. The ligation products were purified and used as template for RNA *in-vitro* transcription. sGluc, secreted *Gaussia* luciferase; 2A, foot-and-mouth disease virus (FMDV) 2A peptide; BSD, blasticidin. **d** Huh7 cells were co-transfected with *in-vitro* transcribed replicon RNA (WT or the indicated mutants) and an mRNA encoding the SARS-CoV-2 N protein. The luciferase activity in the supernatants was measured at the time points indicated. Medium was changed at 8 hours post-transfection. Data are shown as mean±SEM (n=8). SAA, the NSP12 polymerase active-site mutant. Unpaired Student's t-test was performed between the mutants and wild type (WT) and between the mutants as indicated (statistical significance level: *p value<0.05, **p value<0.01, ***p value<0.001).