1    **Analysis of the mutation dynamics of SARS-CoV-2 reveals the spread history and emergence of**

2    **RBD mutant with lower ACE2 binding affinity**

3    Yong Jia[1,*,†], Gangxu Shen[2,3,*], Stephanie Nguyen[4], Yujuan Zhang[1], Keng-Shiang Huang[2,5], Hsing-Ying Ho[6], Wei-Shio

4    Hor[7], Chih-Hui Yang[5], John B Bruning[4], Chengdao Li[1,8,†], Wei-Lung Wang[3,†]

5    [1]College of Science, Health, Engineering and Education, Murdoch University, Murdoch, WA, 6150, Australia

6    [2]School of Chinese Medicine for Post-Baccalaureate, I-Shou University, Kaohsiung, Taiwan

7    [3]National Changhua University of Education, Changhua, Taiwan

8    [4] Institute of Photonics and Advanced Sensing (IPAS), School of Biological Sciences, University of Adelaide, Adelaide

9    5005, Australia

10   [5]College of Medicine, I-Shou University, Kaohsiung, Taiwan

11   [6]Guo-Yuan Clinic, Taichung, Taiwan

12   [7]TOPO Pharmaceutical Co., Ltd, Taichung, Taiwan

13   [8]Department of Primary Industry and Regional Development, Government of Western Australia, South Perth, WA, 6155,

14   Australia

15   *These authors have contributed equally to the study

16

17   [†]Correspondence author:

18   Dr. Yong Jia: y.jia@murdoch.edu.au

19   Prof. Chengdao Li: c.li@murdoch.edu.au

20   Prof. Wei-Lung Wang: wlwang@cc.ncue.edu.tw

## Summary

Monitoring the mutation dynamics of SARS-CoV-2 is critical for the development of effective approaches to contain the pathogen. By analyzing 106 SARS-CoV-2 and 39 SARS genome sequences, we provided direct genetic evidence that SARS-CoV-2 has a much lower mutation rate than SARS. Minimum Evolution phylogeny analysis revealed the putative original status of SARS-CoV-2 and the early-stage spread history. The discrepant phylogenies for the spike protein and its receptor binding domain proved a previously reported structural rearrangement prior to the emergence of SARS-CoV-2. Despite that we found the spike glycoprotein of SARS-CoV-2 is particularly more conserved, we identified a receptor binding domain mutation that leads to weaker ACE2 binding capability based on in silico simulation, which concerns a SARS-CoV-2 sample collected on 27th January 2020 from India. This represents the first report of a significant SARS-CoV-2 mutant, and requires attention from researchers working on vaccine development around the world.

## Highlights

- Based on the currently available genome sequence data, we provided direct genetic evidence that the SARS-COV-2 genome has a much lower mutation rate and genetic diversity than SARS during the 2002-2003 outbreak.

- The spike (S) protein encoding gene of SARS-COV-2 is found relatively more conserved than other protein-encoding genes, which is a good indication for the ongoing antiviral drug and vaccine development.

- Minimum Evolution phylogeny analysis revealed the putative original status of SARS-CoV-2 and the early-stage spread history.

- We confirmed a previously reported rearrangement in the S protein arrangement of SARS-COV-2, and propose that this rearrangement should have occurred between human SARS-CoV and a bat SARS-CoV, at a time point much earlier before SARS-COV-2 transmission to human.

- We provided first evidence that a mutated SARS-COV-2 with reduced human ACE2 receptor binding affinity have emerged in India based on a sample collected on 27th January 2020.

**Keywords**: Binding affinity, Human ACE2, Mutant, Minimum Evolution, Phylogeny, Receptor binding domain, SARS-CoV-2, Spike protein.

## Introduction

The outbreak of severe acute respiratory syndrome–coronavirus 2 (SARS-CoV-2) has caused an unprecedented pandemic and severe fatality around the world. As of 4th April 2020, the total number of SARS-CoV-2 infection has reached over 1.05 million cases globally, claiming 56,985 lives (Coronavirus disease 2019, Situation Report-15, WHO). Most concerning is that this number is predicted to continue to rise significantly in the coming months. Scientists have been working diligently to understand how the virus spreads and evolves. There is an imminent challenge to develop effective approaches to contain the rapid spread of this pathogen.

In addition to the traditional control methods, such as travel bans and self-isolation, which have clear negative impacts on the economy and disrupt normal social activities, the development of antiviral drugs and vaccines should be the ultimate solution to contain the epidemic and reduce the fatality[1,2]. Similar to other SARS-like CoVs [3,4], SARS-CoV-2 uses its spike (S) protein to bind and invade human cells [5,6]. The S protein and its host receptor are the key targets for drug design and vaccine development [7,8]. Recently, several 3D protein structures of the receptor binding domain (RBD) of SARS-CoV-2 spike protein have been determined [5,6,9]. Elucidation of the structural basis of receptor recognition by SARS-CoV-2 has laid the foundation for future vaccine development [6,9].

Vaccines utilize the human immune system and is specific to the viral-encoded peptides [10]. One of the major concerns for antiviral vaccine development is the constant emergence of new mutations, which may reduce its efficacy in future epidemics [7,10]. A prominent example is Influenza virus in which mutations arise every year, requiring annual immunization [11]. SARS-CoV-2 is a single-stranded RNA virus, whose genome can readily mutate as the virus spreads [12,13]. Interestingly, initial assessment of the first 9 SARS-CoV-2 genome sequences revealed a low level of mutation rate [14]. Several more recent studies also highlighted relatively low genetic diversity and stable genomes for SARS-CoV-2 [15-17], which suggests that only a single vaccine may be required for SARS-CoV-2. However, these results may be based on limited genomic data in the early stage of virus development. It is critical to study and monitor the mutation dynamics of SARS-COV-2 to gain a more accurate understanding of the virus and therefore guide vaccine development.

73  Taking advantage of the increasing amount of genomic data collected around the world, we set to explore the current status

74  of SARS-CoV-2 genomic diversity, assess the mutation rate, and potentially identify the emergence of novel mutations that

75  may require close attention. A total of 106 complete or near complete SARS-CoV-2 genome data covering over 12 countries

76  was downloaded from a public database. The genetic diversity profile and evolutionary rate for each protein-encoding gene

77  was characterized. Phylogenetic analyses in this study revealed clues to the spread history of SARS-CoV-2 in some

78  countries. Most importantly, we identified a SARS-CoV-2 mutation with likely reduced human angiotensin-converting

79  enzyme 2 (ACE2) binding affinity. We confirmed that SARS-CoV-2 has a relatively low mutation rate and suggest that

80  novel mutation with likely varied virulence and different immune characteristics may also emerge.

81

82  **Methods**

83  **Sequence retrieval**

84  The latest sequence data for SARS-CoV-2 and SARS was retrieved from the NCBI public database at

85  https://www.ncbi.nlm.nih.gov/genbank/sars-cov-2-seqs/. The 5'UTR, 3'UTR, and CDS sequences of the reference SARS-

86  CoV-2 genome (NC_045512.2) and the human SARS genome (NC_004718.3) were used to blastn against the available

87  genome data. The homology search targets were restricted to the complete or near-complete genomes for further analyses.

88  **Conservation profiling**

89  The assessment of sequence conservation was performed using the PLOTCON tool from the The European Molecular

90  Biology Open Software Suite at https://www.bioinformatics.nl/cgi-bin/emboss/plotcon. The gene model of SARS-CoV-2

91  was generated using the AnnotationSketch [18] tool based on the genome annotation data downloaded from NCBI database.

92  **Phylogeny construction**

93  Codon-based sequence alignment was performed for the conserved domain sequences (CDS) using the MUSCLE program

94  (limited to 2 iterations for fast alignment of long sequences) [19]. Alignment of the 5'UTR and 3'UTR sequences were

95  performed separately. The obtained alignment files were concatenated for final phylogeny construction. The phylogenetic

96  tree was developed in MEGA7.0 [20] using the Minimum Evolution method with p-distance substitution model, and the

97  Maximum Likelihood method (HKY+G+I, 500 times bootstrap test) for the S protein analyses. Tree annotation was carried

98    out using Figtree software ( http://tree.bio.ed.ac.uk/software/figtree/ ) and cophyloplot from ape 5.0 R package [21].

**Evolutionary rate assessment**

100   The ratio of nonsynonymous mutations ($d_N$) to synonymous mutations ($d_S$) was calculated using codeml in the PAML

101   (version 4.7) package [22]. CDS sequences for each protein encoding gene were filtered to remove redundant identical

102   sequences. Then codon-based CDS sequence alignment was performed using the MUSCLE program, and an individual NJ

103   tree was generated using MEGA7.0 [20] with p-distance model. The obtained sequence alignment and phylogenetic tree files

104   were used as PAML inputs for $d_N$ and $d_S$ calculations.

**Protein structural analyses**

106   3D structure of the SARS-CoV-2 spike glycoprotein in complex with human ACE2 (PDB: 6VW1) has been determined

107   recently [5,9]. The structural model for the receptor binding domain (RBD) was extracted from 6VW1 for comparison analysis

108   with the human SARS structure (PDB: 2AJF) [3], which is in complex with the receptor: human ACE2. Amino acid sequence

109   alignment     of     the     spike     glycoprotein     was     visualized     and     annotated     using     ESPript     3.0     tool

110   ( http://espript.ibcp.fr/ESPript/ESPript/index.php ). Protein hydrophobicity profiles were implemented in PyMOL using

111   the Color_h script (http://www.pymolwiki.org/index.php/Color_h), based on the hydrophobicity scale defined at

112   http://us.expasy.org/tools/pscale/Hphob.Eisenberg.html. All structure visualization was carried out using PyMol (Version

113   2.2.3. Schrodinger, LLC).

**In silico mutagenesis and prediction of change in binding free energy**

115   The crystal structure of SARS-CoV-2 spike glycoprotein in complex with the human receptor, angiotensin-converting

116   enzyme 2 (ACE2) (PDB: 6VW1) was used to generate a model of the R408I SARS-CoV-2 spike glycoprotein mutant using

117   ICM-Pro (Molsoft LLC, La Jolla, CA, USA). The model was subsequently refined through the optimization of geometric

118   restraints, refinement of clashing side-chains and minimization of free energy. Prediction of the binding free energy change

119   of the SARS-CoV-2 spike glycoprotein, wild-type and R408I mutant, and the ACE2 receptor interaction was performed

120   using ICM-Pro. All structure visualization was carried out using PyMol (Version 2.2.3. Schrodinger, LLC).

121

## Results

**Genetic diversity analyses identified a single amino acid mutation in RBD of the spike protein in SARS-CoV-2**

As of 24th March 2020, there are a total of 174 nucleotide sequences for SARS-CoV-2 in the NCBI database. By restricting to the complete or near-complete genomes, 106 sequences from 12 countries were obtained and used for further analyses. This encompasses 54 records from USA, 35 from China, and the rest from other countries: Australia (1), Brazil (2), Finland (1), India (2), Italy (1), Japan (3), Nepal (1), Spain (3), South Korea (1), and Sweden (1).

Based on the gene model of the reference SARS-CoV-2 genome (GeneBank: NC_045512.2), a total of 12 protein-encoding open reading frames (ORFs), plus the 5'UTR and 3'UTR were annotated (**Figure 1A**). Overall, the gene sequences from different samples are highly homologous, sharing > 99.1% identity, apart from the 5'UTR (96.7%) and 3'UTR (98%) (**Table 1**), which are relatively more divergent. Sequence alignment showed that there is no mutation in ORF6, ORF7a, and ORF7b. The genetic diversity profile across the 106 genomes was displayed in **Figure 1A**. A few nucleotide sites within ORF1a, ORF1b, ORF3a, and ORF8 exhibiting high genetic diversity were identified (**Figure 1A**).

The S protein is critical for virus infection and vaccine development. As shown in **Figure 1B**, 12 single amino acid substitutions in 12 genomes were identified for the spike glycoprotein, only one (R408I) of which occurs in the receptor binding domain (RBD). This mutation concerns an accession collected from Kerala State, India on 27th Jan 2020.

To track the occurrence of the R408I mutant, we checked the latest GISAID database (5th Oct 2020) and confirmed that there are a total of 17 SARS-CoV-2 strains containing the R408I mutation (**Table 1**): England (11), Egypt (2), Portugal (1), Switzerland (1), and India (2). We believe that these numbers are still underestimated by the limited sequencing capacity in respective countries. For example, there are only a total of 152 spike protein records for Egypt in the GISAID database. Noteworthy, the latest R408I SARS-CoV-2 samples were collected on 10th Sep 2020 and 1st Sep 2020 from Switzerland and England (**Table 1**), respectively, indicating that this mutant is still actively spreading.
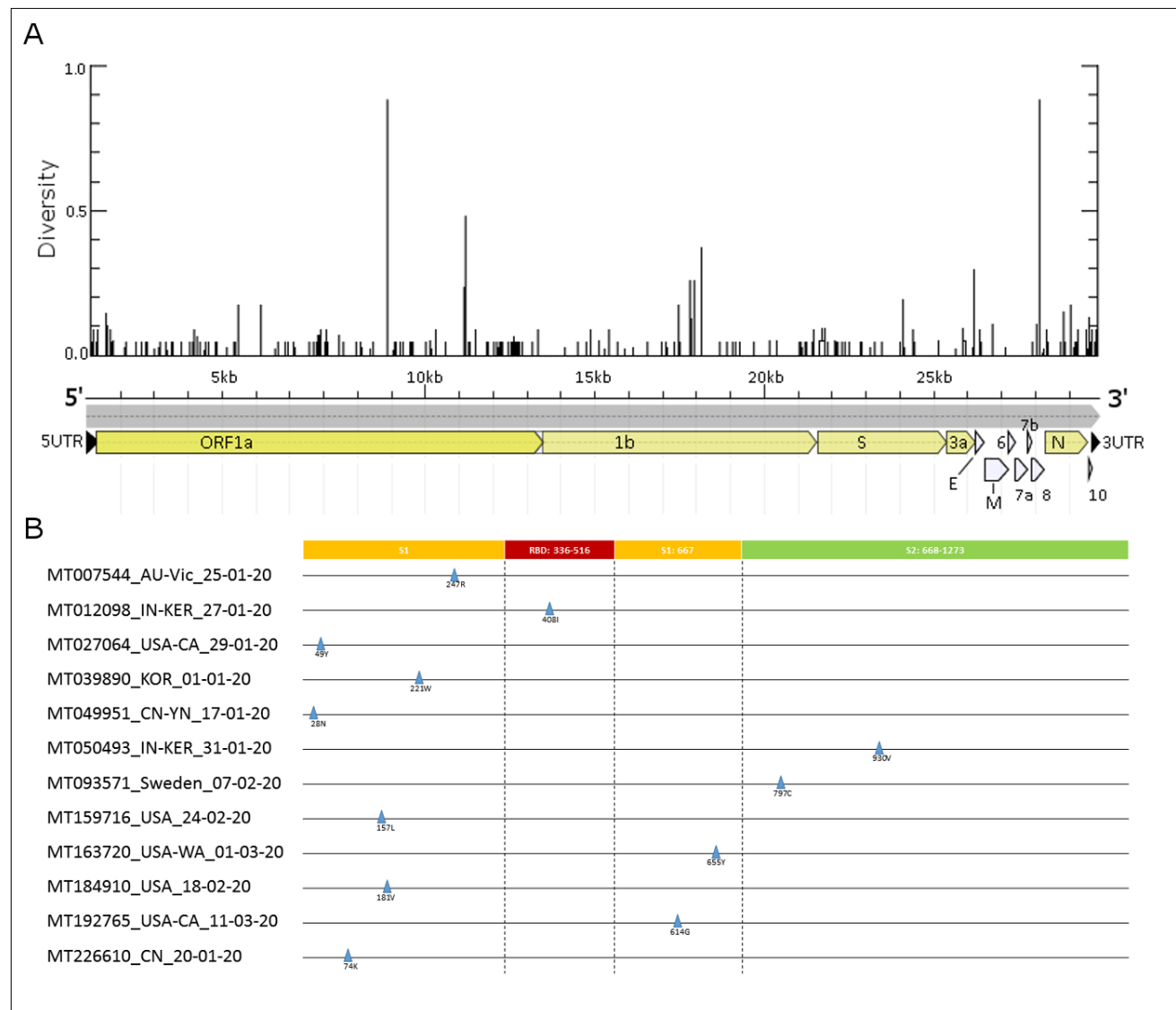
143

**Figure 1. Genetic diversity profile of SARS-CoV-2 genomes and amino acid mutations in the spike glycoprotein. A**) Pair-wise genetic distance for each nucleotide site calculated from the 106 SARS-CoV-2 genomes. Gene model is based on the reference genome (GeneBank: NC_045512.2). **B**) Identification of amino acid mutations in the spike glycoprotein. Sequences were named as: Accession name_country_ sample collection time (AU: Australia; IN: India; USA: United States; KOR: South Korea; CN: China; Sweden: Sweden.) Amino acid numbering according to the reference spike protein (Accession ID: YP_009724390.1).

150 **Table 1. List of SARS-CoV-2 strains containing the R408I mutation**. GISAID and NCBI databases (as of 5th Oct 2020) were searched for strains

151 containing the R408I mutation.

| No | Sample name | Accession ID | Collection date | Location |
|---|---|---|---|---|
| 1 | hCoV-19/Switzerland/BE-ETHZ-280249/2020 | EPI_ISL_560454 | 10/09/2020 | Switzerland / Bern |
| 2 | hCoV-19/England/MILK-9AA096/2020 | EPI_ISL_550882 | 1/09/2020 | England |
| 3 | hCoV-19/England/OXON-B1C55/2020 | EPI_ISL_479081 | 30/06/2020 (Submission) | England |
| 4 | hCoV-19/Egypt/CUNCI-HGC3I013/2020 (MT627395 NCBI database) | EPI_ISL_479694 (MT627395.1) | 2/06/2020 | Egypt |
| 5 | hCoV-19/England/NORT-29DB8D/2020 | EPI_ISL_484309 | 7/05/2020 | England |
| 6 | hCoV-19/England/NORT-29E005/2020 | EPI_ISL_488118 | 6/05/2020 | England |
| 7 | hCoV-19/England/NORW-E7A01/2020 | EPI_ISL_449088 | 4/05/2020 | England |
| 8 | hCoV-19/England/NORT-29D437/2020 | EPI_ISL_499803 | 3/05/2020 | England |
| 9 | hCoV-19/Egypt/CUNCI-HGC013/2020 (MT510693 NCBI database) | EPI_ISL_468047 (MT510693.1) | 2/05/2020 | Egypt |
| 10 | hCoV-19/England/NORW-EE30F/2020 | EPI_ISL_490529 | 24/04/2020 | England |
| 11 | hCoV-19/England/NORT-29C84B/2020 | EPI_ISL_488132 | 23/04/2020 | England |
| 12 | hCoV-19/England/PORT-2D11E5/2020 | EPI_ISL_475338 | 21/04/2020 | England |
| 13 | hCoV-19/England/OXON-B07DD/2020 | EPI_ISL_478909 | 20/04/2020 | England |
| 14 | hCoV-19/Portugal/PT0716/2020 | EPI_ISL_510975 | 24/03/2020 | Portugal |
| 15 | hCoV-19/England/CAMB-74A09/2020 | EPI_ISL_425342 | 18/03/2020 | England |
| 16 | MT050491 (NCBI database) | MT050491.1 | 30/1/2020 | India / Kerala |
| 17 | hCoV-19/India/MH-1-27/2020 (MT012098, NCBI database) | EPI_ISL_413522 (MT012098.1) | 27/01/2020 | India / Kerala |

152

### SARS-CoV-2 displayed a much lower mutation rate than SARS-CoV, with a highly conserved S gene

154 To assess the mutation rate and genetic diversity of SARS-CoV-2, the ratio of nonsynonymous mutations ($d_N$) and

155 synonymous mutations ($d_S$) was calculated for each protein-encoding ORF based on the 106 SARS-CoV-2 and 39 SARS

156 genomes. For SARS-CoV-2, the highest $d_N$ was observed for ORF8 (0.0111), followed by ORF1a (0.0081), ORF9 (0.0079),

157 and ORF4 (0.0063) (**Table 2**), indicating these genes may be more likely to accumulate nonsynonymous mutations. In

158 contrast, ORF1b (0.0029), S gene (0.0040) encoding the spike protein, and ORF5 (0.0023) are relatively more conserved

159 in terms of nonsynonymous mutation. Noteworthy, ORF6, ORF7ab and ORF10 are strictly conserved with no

160 nonsynonymous mutation. Compared to SARS-CoV-2, SARS displayed higher mutation rates for all of the ORFs in the

161 virus genome (Table 1), suggesting an overall higher level of genetic diversity and mutation rate. In particular, the $d_N$ and

162 $d_S$ values for the S gene in SARS-CoV is approximately 12 and 7 times higher than that for SARS-CoV-2, respectively. In

163 contrast, the mutation rate differences for ORF1a and ORF1b between SARS-CoV-2 and SARS are relatively milder,

164 varying from 1.5 times to 4.8 times only. In contrast to SARS-CoV-2, which has strictly conserved ORF6, ORF7a, and

165 ORF7b, SARS displayed mutation rates at different levels. Notably, the $d_S$ for ORF10 are comparable between the two

166 genomes at 0.0326 and 0.0341, respectively.

167 **Table 2. Mutation rate analysis on SARS-CoV-2 genes.** Gene model is according to the SARS-CoV-2 reference genome (GeneBank: NC_045512.2).

168 S: spike glycoprotein. "Pair-wise identity" indicate the minimum pair-wise sequence identity among the 106 genomes. $d_N$: nonsynonymous mutation; $d_S$:

169 synonymous mutations. "--": not applicable.

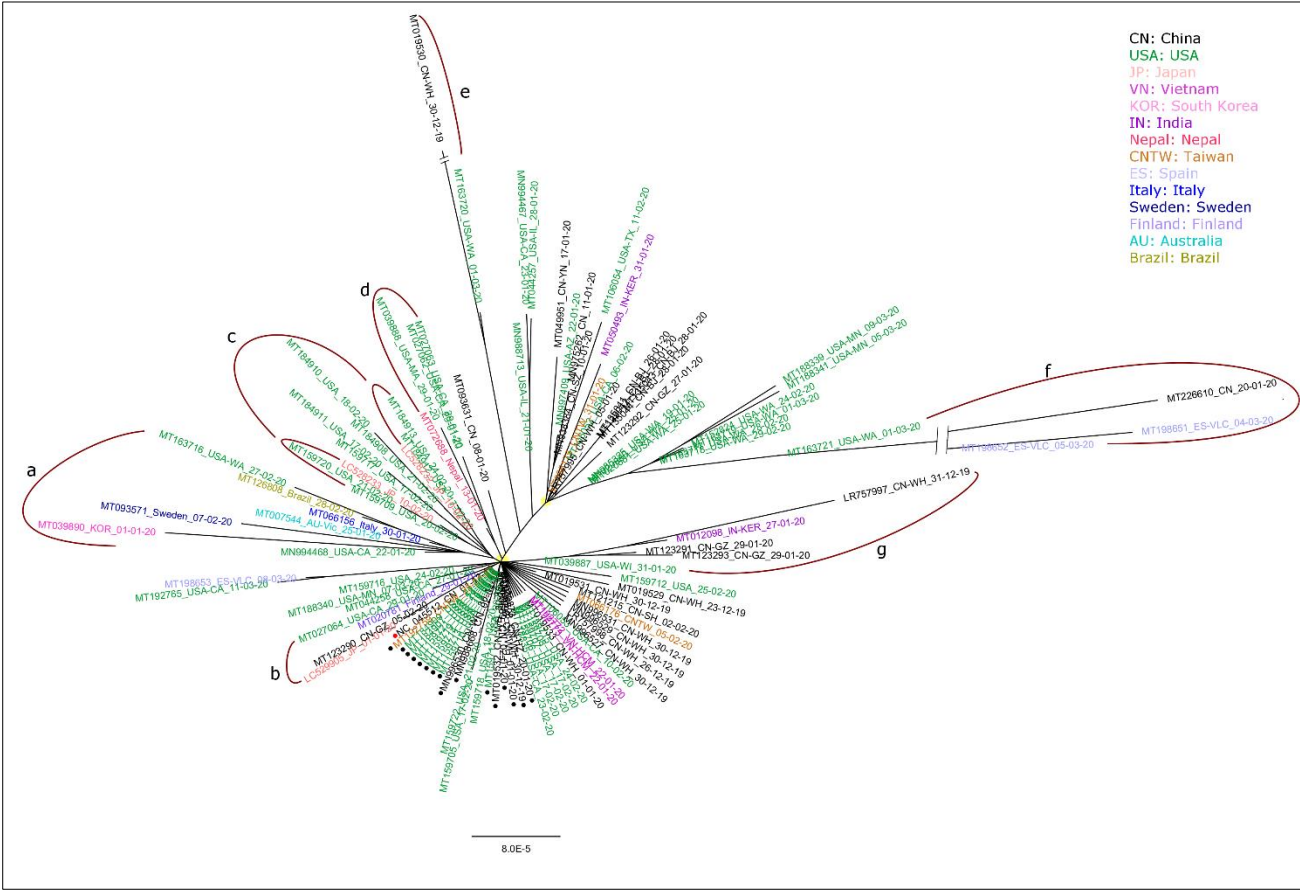| Gene name | | 5'UTR | 1a | 1b | S | ORF3a | ORF4_E | ORF5_M | ORF6 | ORF7a | ORF7b | ORF8 | ORF9 | ORF10 | 3'UTR |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Length (bp) | | 211 | 13218 | 8088 | 3822 | 828 | 228 | 669 | 186 | 336 | 132 | 366 | 1260 | 117 | 152 |
| Pair-wise identity | | 96.7% | 99.8% | 99.9% | 99.9% | 99.6% | 99.1% | 99.7% | 100% | 100% | 100% | 99.5% | 99.7% | 99.1% | 98% |
| $d_N$ | SARS-CoV-2 | -- | 0.0081 | 0.0029 | 0.0040 | 0.0074 | 0.0063 | 0.0023 | 0 | 0 | 0 | 0.0111 | 0.0079 | 0 | -- |
| | SARS | -- | 0.0119 | 0.0077 | 0.0532 | 0.0331 | 0.0338 | 0.023 | 0.3031 | 0.0040 | 0.5339 | 0.0287 | 0.0197 | 0.0135 | -- |
| $d_S$ | SARS-CoV-2 | -- | 0.0041 | 0.0083 | 0.0055 | 0 | 0.0611 | 0.0046 | 0 | 0 | 0 | 0 | 0.0172 | 0.0326 | -- |
| | SARS | -- | 0.0196 | 0.0326 | 0.0442 | 0.0248 | 0.0146 | 0.0928 | 0.0202 | 0.0183 | 0.0005 | 0.0566 | 0.9552 | 0.0341 | -- |

170

171 **Phylogeny analysis revealed the original status of SARS-CoV-2 and its spread history**

172 To trace the potential spread history of SARS-CoV-2 across the world, an unrooted Minimum Evolution (ME) tree of the

173 106 genomes was developed based on whole-genome sequence alignment. The clustering pattern of the ME phylogeny

174 (**Figure 2**) shed light on how the virus may have spread at the early stage. At the center of the ME tree, a number of virus

175 accessions collected from China (including the reference genome NC_045512.2) and USA have the shortest branch

176 (marked by red and black dots), thus may indicate the original status of SARS-CoV-2. The radial pattern, instead of

177 clustering together, of these accessions and other accessions derived from the tree center (highlighted in yellow color) with

178 longer branches, implies the independent mutations occurring during the virus spread (**Figure 2**). However, the longer

179 branch may not be always associated with a longer evolution time, as some accessions collected in December 2019 have

180 equal or even longer branch than those collected in January and February 2020.

181 Due to the data availability, virus accessions collected from China and USA are dominant in the ME tree and constantly

182 group with accessions from other countries. Overall, the target SARS-CoV-2 genomes tend to separate into two major

183 clusters (highlighted in yellow dots, **Figure 2**), suggesting these SARS-CoV-2 may have originated from two major spread

184 sources. Of particular interest is the observation of several phylogenetic clades encompassing samples collected from more

185 than one country, which may provide clues to track the spread history of SARS-CoV-2 in these countries. For example, a

186 notable clade (clade a) containing accessions collected from USA, Brazil, Italy, Australia, Sweden, and South Korea was

187 identified. The only Brazil accession (MT126808.1) in this study is found to be clustered with one accession from USA

188    (MT163716.1) with strong support. Whilst the virus accessions from China are prevalent in the ME tree, it is intriguing

189    that no correlated accession from China is found in this clade. An additional clade including accessions collected from

190    China, USA and Finland were found together (clade b). In another notable clade (clade c), 2 of the 3 accessions

191    (LC528232.1 and LC528233.1) collected from the cruise ship in Japan were grouped with several accessions from USA.

192    Two accessions (MT198651.1 and MT198652.1) collected in March 2020 from Spain were grouped (clade f) with one

193    accession collected in January 2020 from China. The additional Spain accession (MT198653.1) was clustered with one

194    from USA (MT192765.1). One India accession (MT012098.1) was found together (clade g) with an accession from Wuhan,

195    China, collected in December 2019. Interestingly, the single Nepal accession (MT072688.1) seems to be closely related

196    (clade d) to several accessions from USA.



197

198    **Figure 2. Phylogeny clustering analyses of the 106 SARS-CoV-2 genomes**. Results were based on whole genome sequence alignment using Minimum

199    Evolution method. Each accession was named in the "accession ID, country, sample collection time" format. Samples collected from different countries

200    were highlighted in different colors. Red dots indicated the reference SARS-CoV-2 genome (GeneBank: NC_045512.2), which together with black dots

201    indicated the original status of SARS-CoV-2 (branch length = 0). The putative two types of SARS-CoV-2 were highlighted in yellow shades. Notable

202    clades containing sequences from more than one country were highlighted in curved line (magenta).

203 **Spike protein of SARS-CoV-2 has undergone a structural rearrangement**

204 The spike glycoprotein is critical for virus infection. Recent study suggested that the S protein in SARS-CoV-2 may have

205 undergone a structural rearrangement[13]. To investigate this hypothesis, two separate phylogenies were developed based on

206 the full-S and RBD sequences, respectively. Human SARS-CoV-2, MERS, and SARS-CoV reference sequences and their

207 close coronavirus homologues identified from various animal hosts were included for the phylogenetic analyses. Overall,

208 the two phylogenies displayed similar clustering patterns, separating into three major clades (**Figure 3)**. SARS-CoV-2 was

209 identified in the same major clade and was clustered most closely with two bat SARS CoVs (highlighted in purple and

210 green colors, **Figure 3)** and the human SARS-CoV (orange color, **Figure 3**). In both phylogenies, SARS-CoV-2 is most

211 closely related to bat_CoV_RaTG13, suggesting SARS-CoV-2 may have originated from bat. However, the evolutionary

212 positions of human SARS-CoV and bat-SL-CoVZ45 were swapped between the full-S and RBD-only phylogenies. In the

213 full-S phylogeny, bat-SL-CoVZ45 is relatively more similar to human SARS-CoV-2, whilst human SARS-CoV is closer

214 to SARS-CoV-2 than bat-SL-CoVZ45. Taken together, these results suggested that the RBD of SARS-CoV-2 is more likely

215 originated from human SARS-CoV, whilst the remaining part of the S protein in SARS-CoV-2 may have originated from

216 bat-SL-CoVZ45, supporting the potential structural rearrangement of S protein in SARS-CoV-2. bat_CoV_RaTG13 is

217 similar to SARS-CoV-2, indicating the proposed structural rearrangement may have occurred in bat first before its
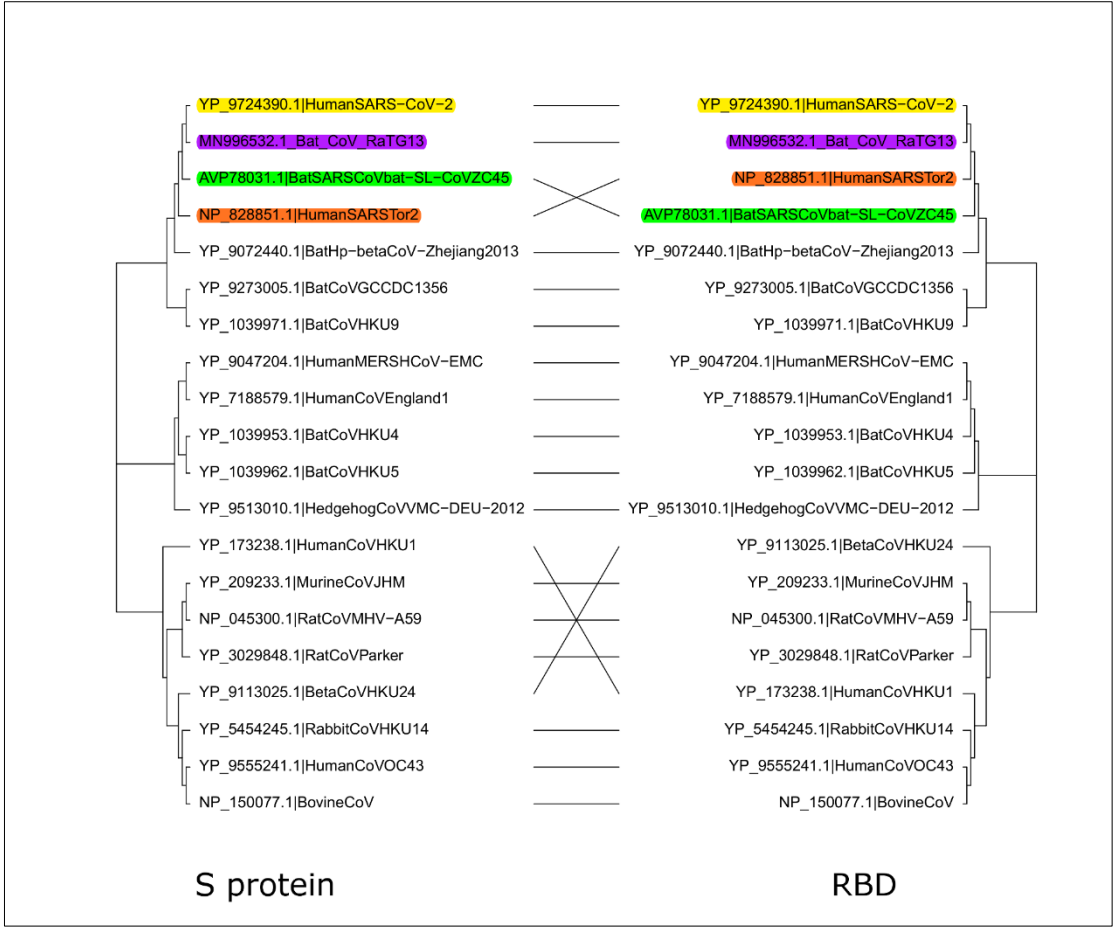
218 transmission to human.

**Figure 3. Displays the phylogeny discrepancy of the full-S and RBD sequences.** Maximum Likelihood phylogenies based on the full-S protein (left) and RBD (right) sequences of SARS-like CoVs. Taxa names were in the "Accession Id, host organism, sample name" format. Human SARS-CoV-2 and its close relatives were highlighted in different colors.

**A single amino acid mutation in RBD results in reduced receptor binding affinity on human ACE2**

The RBD of virus S protein binds to a receptor in host cells, and is responsible for the first step of CoV infection [3]. Thus, amino acid mutation to RBD may have significant impact on receptor binding and vaccine development. The 3D structure of the spike protein RBD of SARS-CoV-2 (PDB: 6VW1) has recently been determined in complex with human ACE2 receptor [6]. One of the 12 spike protein mutations identified above (**Figure 1B**) was located in the RBD region (R408I). Further data screening against the latest GISAID and NCBI database (5th Oct) revealed a total of 17 strains from five countries containing the R408I mutation (**Table 1**). Sequence alignment showed that 408R is strictly conserved in SARS-CoV-2, SARS-CoV and bat SARS-like CoV (**Figure 4A**). Based on the determined CoV2_RBD-ACE2 complex structure, 408R is located at the interface between RBD and ACE2, but is positioned relatively far away from ACE2 and thus does not have direct interaction with ACE2 (**Figure 4B**). However, the determined RBD0-ACE2 structure showed that 408R

234    forms a hydrogen bond (3.3 Å in length) with the glycan attached to 90N from ACE2 (**Figure 4C**) [6]. The hydrogen bond

235    may have contributed to the exceptionally higher ACE2 binding affinity. The arginine residue is also conserved in human

236    SARS-CoV (corresponding to 395R in PDB: 2AJF), but is positioned relatively distant (6.1 Å) from the glycan bound to

237    90N from ACE2 (**Figure S1**). Interestingly, the 408R-glycan hydrogen bond appears to be disrupted by the R408I mutation

238    in one SARS-CoV-2 accession (GeneBank ID: MT012098.1) (**Figure 4D**), collected from India on 27th Jan 2020. *In silico*

239    calculations indicatethat the R408I mutation increased the binding free energy by 0.93 kcal/mol, predicting a modest

240    decrease in ACE2-binding affinity. In contrast to the electrically charged and highly hydrophilic arginine residue, the

241    mutated isoleucine residue has a highly hydrophobic side chain with no hydrogen-bond potential (**Figure 4E**). In summary,

242    the R408I mutation identified from the SARS-CoV-2 strain in India represents a SARS-CoV-2 mutant with likely lower
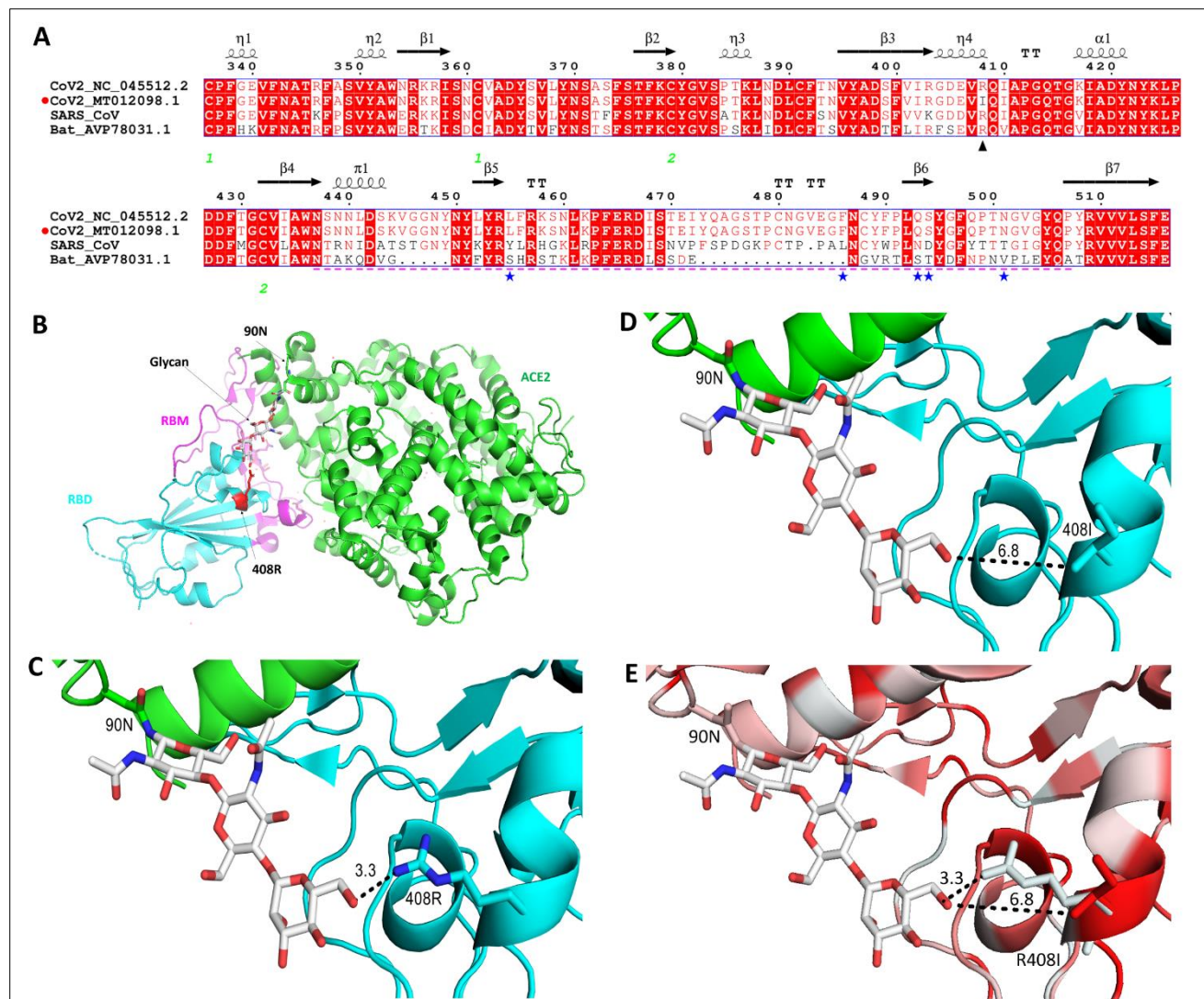
243    ACE2 binding affinity.

244

**Figure 4. Sequence alignment and protein structural analyses of the mutation in RBD of SARS-CoV-2. A**) Sequence alignment of RBDs. ▲: R408I mutation; --- : receptor binding motif (RBM). ★ : RBD-interacting sites. **B**) Overall position of the identified mutation relative to: RBD (cyan), ACE2 (green) with RBM (pink) and Glycan (grey). **C,D**) Display the disrupted hydrogen bond by the R408I mutation. "---": distance in Å. **E**) Hydrophobic profile changes due to R408I mutation, with red and white colors representing the highest hydrophobicity and the lowest hydrophobicity, respectively. All amino acid number according to the S protein of SARS-CoV-2 (NC_045512.2) and human ACE2, respectively.

## Discussions

Based on the currently available genome sequence data, our results showed that the mutation rate of SARS-CoV-2 is much lower than that for SARS, which caused the 2002-2003 outbreak. Our study is the first to provide a direct gene-based quantitative comparison between SARS-COV-2 and SARS. Among the different genetic regions of SARS-CoV-2 genomes, we found that the spike protein (S) is more conserved that other genetic regions such as ORF1, ORF8, and N, which encode nonstructural polyprotein, virus accessory protein, and the nucleocapsid protein, respectively. A relatively stable spike

258 protein region of SARS-CoV-2 is a good indication for the epidemic control, as less mutation raises the hope of the rapid

259 development of a vaccine and antiviral drugs. Our results are consistent with several recent genetic variance analyses on

260 SARS-CoV-2 [15,23-25], which suggested the SARS-CoV-2 genomes are highly homogeneous. Furthermore, based on the

261 latest genomic data for SARS-CoV-2, molecular geneticists monitoring the virus development also suggested that the

262 mutation rate of SARS-CoV-2 maintains at a low level [17,26,27]. Whilst it is generally safe to say that SARS-CoV-2 tends to

263 mutate at a low rate, as the virus continues to spread rapidly around the world, and more genomic data is accumulated, the

264 evolution and mutation dynamics of SARS-CoV-2 still need to be monitored closely.

265

266 One critical aim of our study is to identify the original status of SARS-CoV-2 before its wide transmission across different

267 countries. Due to the short time space of sample collection and a relatively low mutation rate for SARS-CoV-2, we believe

268 that Minimum Evolution phylogeny (a parsimony method) may outperform other phylogenetic methods to achieve this

269 aim. Similarly, Peter et al. [28] also adopted a parsimony phylogeny (Maximum Parsimony) to trace the spread history of

270 SARS-CoV-2 in the early stage of the pandemic. Minimum Evolution and Maximum Parsimony are similar phylogeny

271 methods (both using the parsimony sites detected in the sequence alignment) trying to minimize the total number of

272 substitutions in the phylogenetic tree. In our analysis, the earliest few reported SARS-CoV-2 accessions collected from

273 Wuhan China were identified at the center of the phylogenetic tree with the shortest branch. Interestingly, several virus

274 genomes from USA were found to be identical to these putative original versions of the virus from Wuhan. According to

275 public media, the outbreak of SARS-CoV-2 in USA occurred relatively later than other countries. One possible explanation

276 for this observation is that, the spread of SARS-CoV-2 in USA might start much earlier than previously thought or reported.

277 Since a dominant proportion of the samples in this study were collected from China and USA, we observed a significantly

278 higher level of genetic diversity from these two countries. Most SARS-CoV-2 accessions from the other countries can find

279 their closely related sisters from either China or USA. This data bias, on the other hand, may give us an advantage to trace

280 the spread history of SARS-CoV-2 in different countries. This suggestion is reliable because all samples investigated in

281 this study were collected at the early stage of the pandemic, which may avoid the potential data noise caused by recent

282 published genomes of complex spread background. One notable finding in our phylogenetic tree is that, the singleton

283 SARS-CoV-2 accessions collected from Australia, Brazil, South Korea, Italy and Sweden were clustered together with two

284  USA samples but without a Chinese version, suggesting that these infection cases may be somehow related. In addition,

285  one of the three samples collected from the cruise ship stranded in Japan was found to be closely related to a sample

286  collected from Guangzhou, China, whilst the other two were grouped with several cases from USA. Noteworthy, our

287  phylogeny seems to support the presence of two major types of SARS-CoV-2 in the target samples, suggesting the potential

288  existence of two spread sources. Interestingly, this speculation is corroborated by an independent clustering analyses using

289  a different phylogeny method [23].

290

291  Until now, the origin of SARS-CoV-2, and how it has been transmitted to humans remains largely a mystery. Early genomic

292  data indicated that human SARS-CoV-2 is an enveloped, positive-sense, and single-stranded RNA virus in the subgenus

293  *Sarbecovirus* of the genus *Betacoronavirus* [13,14]. Evolutionarily, SARS-CoV-2 is most closely related to bat SARS-like

294  CoV (88% genome sequence identity) and human SARS CoV (79%), the latter of which caused a global pandemic in 2003

295  [13]. Based on the strong genome sequence identity between SARS-CoV-2 and bat SARS-like COVs, it was initially

296  speculated that SARS-CoV-2 may have originated from bat [14,29]. However, a more recent study proposed that pangolin may

297  be the most likely reservoir hosts due to the identification of closely related SARS-COVs from this species as well [30]. Both

298  animals can harbor coronaviruses related to SARS-CoV-2. However, direct evidence of the transmission of SARS-CoV-2

299  from either bat or pangolin to human is still missing.

300

301  Prior to this study, several publications have suggested that SARS-CoV-2 may have originated from the genome

302  recombination of SARS-like CoVs from different animal hosts, as evidenced by the discrepant clustering patterns for the

303  phylogenies using different genetic regions. Lu [13] first observed that the RBD of S protein in SARS-CoV-2 is more closely

304  related to human SARS-CoV, whilst the other part of its genome is more similar to bat SARS-CoV. Later Lam [30] identified

305  a bat CoV_RaTG13 and several pangolin SARS-CoVs that are consistently closer to SARS-CoV-2 than human SARS-

306  CoV in either full-S protein or RBD. By combining the data from these two studies, our study confirmed the observations

307  reported in both studies, and further determined that the S protein recombination actually happened between human SARS-

308  CoV and a bat SARS-CoV, much earlier before its transmission to human, with the newly identified bat SARS-CoV-

309  RaTG13 as an intermediate.

310

311  The RBD of S protein binds to a receptor in host cells and is responsible for the first step of CoV infection. The receptor

312  binding affinity of RBD directly affects virus transmission rate. Thus, it has been the major target for antiviral vaccine and

313  therapeutic development such as SARS [8]. At the time of first completion in late March 2020, this study was the first to

314  report the identification of the R408I mutation in the RBD of S protein in SARS-CoV-2. Since then, the R408I mutant has

315  attracted research attention from a significant number of researchers. Both computational and experimental studies have

316  been performed to further investigate its molecular characteristics and potential immune effects [31-38]. In addition,

317  commercial synthesis of the R408I recombinant RBD protein has been offered by serval companies (Acro Biosystems,

318  Creative Dianostics, SinoBiological, and Creative Biolabs) for immuno-binding and diagnostic testing. Noteworthy, Yan

319  et al. [31] showed that three of the four RBD neutralizing antibodies tested could not bind the R408I mutant, whereas other

320  mutants displayed strong binding interaction with all the neutralizing antibodies tested. The authors stated that 408R played

321  an essential role for SARS-CoV-2 RBD antibody binding and the R408I could abolish the antibody binding interaction [31].

322  In addition, Zhe et al. [39] also suggested that R408I constitute the RBD epitope residues. These observations contrast an

323  early stage study [17] which did not notice the R408I mutation and predicted that a single vaccine may be sufficient for all

324  circulating SARS-CoV-2 variant. Based on the determined RBD-hACE2 protein structure (PDB: 6VW1) [6], we found that

325  408R residue can establish an indirect receptor interaction via a glycan attached to human ACE2. This residue was found

326  to be conserved in SARS and MERS as well. Interestingly, the arginine residue (corresponding to 395R in SARS) has also

327  been shown to be involved in receptor interaction in SARS. In this study, we were the first to show that the R408I mutation

328  in SARS-CoV-2 is likely to cause a reduced binding affinity to human ACE2 receptor. Our result has been corroborated in

329  several independent studies later on [32-34]. Although the S protein gene seems to be more conserved than the other protein-

330  encoding genes in the SARS-CoV-2 genome, our study provides direct evidence that a mutated version of SARS-CoV-2 S

331  protein with varied transmission rate may have already emerged. Furthermore, we confirmed that, as of 5th Oct 2020, a

332  total of 17 SARS-CoV-2 strains containing the R408I mutation were present in the GISAID and NCBI databases, with the

333  latest R408I mutants collected on 10th Sep 2020 and 1st Sep 2020 from Switzerland and England, respectively. These results

334  suggest that the R408I mutant may spread across to different countries since its first emergence from India and is still

335  actively spreading in different regions. Benson et al.[40] recently reported that R408I accounts for ~2% of infection cases in

336  Africa. We believe that the number of identified R408I mutants are still underestimated, given the limited sequencing

337  capability in respective countries. Based on the close relationship of SARS-CoV-2 to SARS, current vaccine and drug

338  development for SARS-CoV-2 has also focused on the S protein and its human binding receptor ACE2 [7,41]. Considering

339  the significantly varied antibody binding profile for R408I, we propose that this mutant still requires significant attention

340  from doctors and scientists around the world during the development of SARS-CoV-2 therapeutic solutions. One suggestion

341  for the next step of therapeutic development is to focus on the identification of potential human ACE2 receptor blocker, as

342  suggested in a recent commentary [7]. This approach will avoid the above-mentioned challenge faced by vaccine

343  development.

## Acknowledgement

348

349

## Author Contribution

351  WLW, CL and YJ conceived the study. YJ, GS, SN, JBB, YZ, KSH, HYH, WSH, CHY performed data analyses.

352  YJ,GS&SN wrote the manuscript. All authors have read the manuscript.

## Conflict of interest

354  The authors declare no conflict of interest.
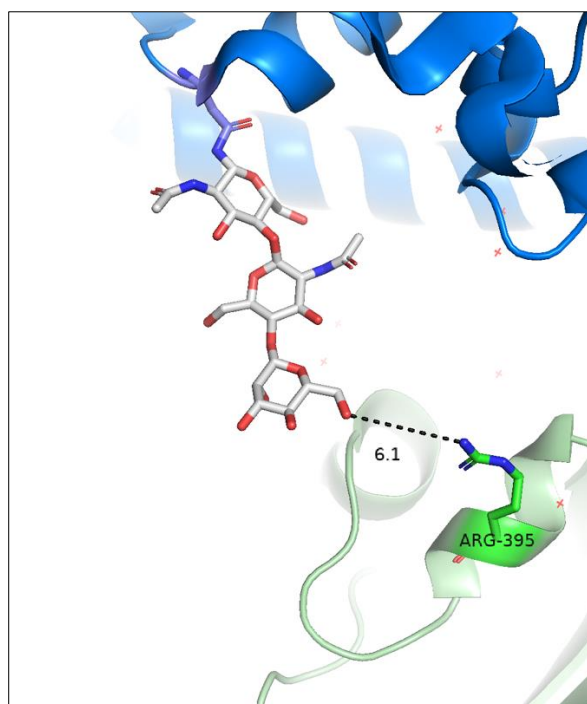
## Supplementary figure



**Figure S1. Displays the position of 395R in human SARS-CoV (PDB: 2AJF).** Dash line indicates the measured distance in Å.

## Data availability

The Genebank ID list of 106 SARS-CoV-2 and 39 SARS genomes used in this study is available at

https://figshare.com/s/3d3c24ef05084b534b4c

## References:

1.    Zhang L, Liu YH. Potential interventions for novel coronavirus in China: A systematic review. *J Med Virol* 2020; **92**(5): 479-90.

2.    Lu S. Timely development of vaccines against SARS-CoV-2. *Emerg Microbes Infec* 2020; **9**(1): 542-4.

3.    Li F, Li WH, Farzan M, Harrison SC. Structure of SARS coronavirus spike receptor-binding domain complexed with receptor. *Science* 2005; **309**(5742): 1864-8.

4.    Li F. Evidence for a Common Evolutionary Origin of Coronavirus Spike Protein Receptor-Binding Subunits. *J Virol* 2012; **86**(5): 2856-8.

5.    Wrapp D, Wang NS, Corbett KS, et al. Cryo-EM structure of the 2019-nCoV spike in the prefusion conformation. *Science* 2020; **367**(6483): 1260-+.

6.    Shang J, Ye G, Shi K, et al. Structural basis of receptor recognition by SARS-CoV-2. *Nature* 2020.

7.    Gurwitz D. Angiotensin receptor blockers as tentative SARS-CoV-2 therapeutics. *Drug development research* 2020.

8.    Du LY, He YX, Zhou YS, Liu SW, Zheng BJ, Jiang SB. The spike protein of SARS-CoV - a target for vaccine and

378    therapeutic development. *Nat Rev Microbiol* 2009; **7**(3): 226-36.

379    9.    Yan R, Zhang Y, Li Y, Xia L, Guo Y, Zhou Q. Structural basis for the recognition of SARS-CoV-2 by full-length
380    human ACE2. *Science* 2020; **367**(6485): 1444-8.

381    10.    Correia B, Bates J, Loomis R, et al. Proof of principle for epitope-focused vaccine design. *Protein Sci* 2015;
382    **24**: 181-4.

383    11.    Huckriede A, Bungener L, Daemen T, Wilschut J. Influenza Virosomes in Vaccine Development. *Methods*
384    *in Enzymology* 2003; **373**: 74-91.

385    12.    Denison MR, Graham RL, Donaldson EF, Eckerle LD, Baric RS. Coronaviruses An RNA proofreading machine
386    regulates replication fidelity and diversity. *Rna Biol* 2011; **8**(2): 270-9.

387    13.    Lu RJ, Zhao X, Li J, et al. Genomic characterisation and epidemiology of 2019 novel coronavirus:
388    implications for virus origins and receptor binding. *Lancet* 2020; **395**(10224): 565-74.

389    14.    Xu XT, Chen P, Wang JF, et al. Evolution of the novel coronavirus from the ongoing Wuhan outbreak and
390    modeling of its spike protein for risk of human transmission. *Sci China Life Sci* 2020; **63**(3): 457-60.

391    15.    Ceraolo C, Giorgi FM. Genomic variance of the 2019-nCoV coronavirus. *J Med Virol* 2020.

392    16.    Tang X. On the origin and continuing evolution of SARS-CoV-2. *Microbiology* 2020.

393    17.    Dearlove B, Lewitus E, Bai H, et al. A SARS-CoV-2 vaccine candidate would likely match all currently
394    circulating variants. *P Natl Acad Sci USA* 2020; **117**(38): 23652-62.

395    18.    Steinbiss S, Gremme G, Schrfer C, Mader M, Kurtz S. AnnotationSketch: a genome annotation drawing
396    library. *Bioinformatics* 2009; **25**(4): 533-4.

397    19.    Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids*
398    *Res* 2004; **32**(5): 1792-7.

399    20.    Kumar S, Stecher G, Tamura K. MEGA7: molecular evolutionary genetics analysis version 7.0 for bigger
400    datasets. *Mol Biol Evol* 2016; **33**(7): 1870-4.

401    21.    Paradis E, Schliep K. ape 5.0: an environment for modern phylogenetics and evolutionary analyses in R.
402    *Bioinformatics* 2019; **35**(3): 526-8.

403    22.    Yang ZH. PAML 4: Phylogenetic analysis by maximum likelihood. *Mol Biol Evol* 2007; **24**(8): 1586-91.

404    23.    Tang X, Wu C, Li X, et al. On the origin and continuing evolution of SARS-CoV-2. *National Science Review*
405    2020.

406    24.    Pachetti M, Marini B, Benedetti F, et al. Emerging SARS-CoV-2 mutation hot spots include a novel RNA-
407    dependent-RNA polymerase variant. *J Transl Med* 2020; **18**(1).

408    25.    Kaushal N, Gupta Y, Goyal M, Khaiboullina SF, Baranwal M, Verma SC. Mutational Frequencies of SARS-
409    CoV-2 Genome during the Beginning Months of the Outbreak in USA. *Pathogens* 2020; **9**(7).

410    26.    Mercatelli D, Giorgi FM. Geographic and Genomic Distribution of SARS-CoV-2 Mutations. *Front Microbiol*
411    2020; **11**.

412    27.    Jaroszewski L, Iyer M, Alisoltani A, Sedova M, Godzik A. The interplay of SARS-CoV-2 evolution and
413    constraints imposed by the structure and functionality of its proteins. *bioRxiv* 2020.

414    28.    Forster P, Forster L, Renfrew C, Forster M. Phylogenetic network analysis of SARS-CoV-2 genomes. *P Natl*
415    *Acad Sci USA* 2020; **117**(17): 9241-3.

416    29.    Zhou P, Yang XL, Wang XG, et al. A pneumonia outbreak associated with a new coronavirus of probable
417    bat origin. *Nature* 2020; **579**(7798): 270-+.
418    30.    Lam TT-Y, Shum MH-H, Zhu H-C, et al. Identifying SARS-CoV-2 related coronaviruses in Malayan pangolins.
419    *Nature* 2020; **583**: 282-5.
420    31.    Lou Y, Zhao W, Wei H, et al. Cross-neutralization antibodies against SARS-CoV-2 and RBD mutations from
421    convalescent patient antibody libraries. *bioRxiv* 2020.
422    32.    Khan MI, Khan ZA, Baig MH, et al. Comparative genome analysis of novel coronavirus (SARS-CoV-2) from
423    different geographical locations and the effect of mutations on major target proteins: Anin silicoinsight. *Plos*
424    *One* 2020; **15**(9).
425    33.    Lokmana SM, Md.Rasheduzzaman, Salauddina A, et al. Exploring the genomic and proteomic variations
426    of SARS-CoV-2 spike glycoprotein: A computational biology approach. *Infection, Genetics and Evolution* 2020;
427    **84**.
428    34.    Ahamad S, Kanipakam H, Gupta D. Insights into the structural and dynamical changes of spike
429    glycoprotein mutations associated with SARS-CoV-2 host receptor binding. *J Biomol Struct Dyn* 2020.
430    35.    Li QQ, Wu JJ, Nie JH, et al. The Impact of Mutations in SARS-CoV-2 Spike on Viral Infectivity and
431    Antigenicity. *Cell* 2020; **182**(5): 1284-+.
432    36.    Laha S, Chakraborty J, Das S, Manna SK, Biswas S, Chatterjee R. Characterizations of SARS-CoV-2
433    mutational profile, spike protein stability and viral transmission. *Infection, Genetics and Evolution* 2020; **85**.
434    37.    Chand GB, Banerjee A, Azad GK. Identification of twenty-five mutations in surface glycoprotein (Spike) of
435    SARS-CoV-2 among Indian isolates and their impact on protein dynamics. *Gene Reports* 2020; **21**.
436    38.    Kim SI, Noh J, Kim S, et al. Stereotypic Neutralizing VH Clonotypes Against SARS-CoV-2 RBD in COVID-19
437    Patients and the Healthy Population. *bioRxiv* 2020.
438    39.    Lv Z, Deng Y-Q, Ye Q, et al. Structural basis for neutralization of SARS-CoV-2 and SARS-CoV by a potent
439    therapeutic antibody. *Science* 2020; **369**(6510): 1505-9.
440    40.    Iweriebor BC, Egbule OS, Danso SO, et al. Analysis of SARS-CoV-2 genomes from across Africa reveals
441    potentially clinically relevant mutations. *bioRxiv* 2020.
442    41.    Ahmed SF, Quadeer AA, McKay MR. Preliminary identification of potential vaccine targets for the COVID-
443    19 coronavirus (SARS-CoV-2) based on SARS-CoV immunological studies. *Viruses* 2020; **12**(3): 254.
444