1 **Phylogeny and evolution of the SARS-CoV-2 spike gene from December 2022 to**

2 **February 2023**

3

4

5

6 Hsiao-Wei Kao[1*]

7

8 [1]Department of Life Sciences, National Chung Hsing University, Taiwan, R.O.C.

9 145 Xingda Road., South District., Taichung City 40227

10 Fax: +886-4-22874740

11 Tel：+886-4-22840416

12

13 *Correspondence: Hsiao-Wei Kao

14 E-mail: hkao@dragon.nchu.edu.tw

15

18

19 **Abstract**

20 **Background:** By the end of 2022, new variants of SARS-CoV-2, such as BQ.1.1.10,

21 BA.4.6.3, XBB, and CH.1.1, emerged with higher fitness than BA.5.

22 **Methods:** The file (spikeprot0304), which contains spike protein sequences, isolates

23 collected before March, 4, 2023, was downloaded from Global Initiative on Sharing

24 All Influenza Data (GISAID). A total of 188 different spike protein sequences were

25 chosen, of which their isolates were collected from December 2022 to February 2023.

26 These sequences did not contain undetermined amino acid X, and each spike protein

27 sequence had at least 100 identical isolate sequences in GISAID. Phylogenetic trees

28 were reconstructed using IQ-TREE and MrBayes softwares. A median-join network

29 was reconstructed using PopART software. Selection analyses were conducted using

30 site model of PAML software.

31 **Results:** The phylogenetic tree of the spike DNA sequences revealed that the majority

32 of variants belonged to three major lineages: BA.2 (BA.1.1.529.2), BA.5

33 (BA.1.1.529.5), and XBB. The median network showed that these lineages had at

34 least six major diversifying centers. The spike DNA sequences of these diversifying

35 centers had the representative accession IDs (EPI_ISL_) of 16040256 (BN.1.2),

36 15970311 (BA.5), 16028739 (BA.5.11), 16028774 (BQ.1), 16027638 (BQ.1.1.23),

37 and 16044705 (XBB.1.5). Selection analyses revealed 26 amino-acid sites under

38    positive selection. These sites included L5, V83, W152, G181, N185, V213, H245,

39    Y248, D253, S255, S256, G257, R346, R408, K444, V445, G446, N450, L452, N460,

40    F486, Q613, Q675, T883, P1162, and V1264.

41    **Conclusion:** The spike proteins of SARS-CoV-2 from December 2022 to February

42    2023 were characterized by a swarm of variants that were evolved from three major

43    lineages: BA.2 (BA.1.1.529.2), BA.5 (BA.1.1.529.5), and XBB. These lineages had at

44    least six diversifying centers. Selection analysis identified 26 amino acid sites were

45    under positive selection. Continued surveillance and research are necessary to monitor

46    the evolution and potential impact of these variants on public health.

47

48    **Keywords:** BA.5, BQ.1, diversifying center, median-join network, PAML, selection,
49    XBB, XBC
50
51

52    [*]Correspondence: hkao@dragon.nchu.edu.tw

53    [1]Department of Life Sciences, National Chung Hsing University, Taiwan, R.O.C.

54    145 Xingda Road., South District., Taichung City 40227

55

**Background**

56

57     On May 5, 2023, the World Health Organization (WHO) declared that COVID-

58     19 is no longer a public health emergency of international concern (PHEIC) due to the

59     decreasing trend in COVID-19 deaths, decline in COVID-19-related hospitalizations

60     and intensive care unit admissions, and the high levels of population immunity to

61     SARS-CoV-2 [1].

62     The Omicron (B.1.1.529) variant was designated as the fifth variant of concern

63     declared by the WHO on November 26, 2021 [2]. A comparison between the B.1.529

64     variant and the Wuhan-Hu-1 genome sequences revealed 53 nucleotide substitutions.

65     Within these substitutions, 30 were nonsynonymous substitutions located in the spike

66     gene [3, 4]. Additionally, there were six amino acid deletions at positions 69, 70, 143,

67     144, 145, and 211. Furthermore, three amino acid insertions (EPE) were observed

68     between positions 214 and 215, relative to the amino acid positions in the Wuhan-Hu-

69     1 spike protein [3, 4].

70     The major lineages that contributed to the pandemic from 2019 to 2022 were

71     Omicron BA.1, BA.2, BA.3, BA.4, and BA.5 [5]. Recently, new variants have

72     emerged, including BQ.1.1.10, BA.4.6.3, XBB, and CH.1.1, which had higher fitness

73     than BA.5 [6-8]. This higher fitness includes evasion of neutralization drugs and

74     convalescent plasma, even those targeting BA.5 breakthrough infections. The immune

4

75    escape mechanism of these new variants is primarily attributed to specific mutations

76    at amino acid sites R346, R356, K444, V445, G446, N450, L452, N460, F486, F490,

77    R493, and S494 within the receptor binding domain of the spike protein. These

78    mutations have been observed in at least five different phylogenetic lineages, which

79    suggests that there has been convergent evolution of the receptor binding domain

80    driven by preexisting SARS-CoV-2 humoral immunity [6-8].

81        In this study, the evolution of the SARS-CoV-2 spike gene between December

82    2022 and February 2023 was investigated. To summarize the major lineages of SARS-

83    CoV-2 and their spike gene evolution during this period, a phylogenetic tree and

84    median-joining network were reconstructed. Furthermore, to identify amino acid sites

85    that were potentially under positive selection and associated with adaptive changes in

86    the spike gene, the nonsynonymous versus synonymous substitution ratio (dn/ds ratio

87    $= \omega$) was calculated. This was done using the site model in the codeml module of the

88    PAML software [9].

89    **Methods**

90    **Data collection and analyses**

91    The file "spikeprot0304" containing spike protein sequences was downloaded from

92    the Global Initiative on Sharing All Influenza Data (GISAID) [10]. To filter the

93    sequences, the following criteria were applied using the Bioedit software [11]: the

94    collection days ranged from December 2022 to February 2023, the sequence lengths

95    ranged from 1259 to 1319 amino acids, and sequences without undetermined amino

96    acid X were included. After filtering, a total of 369,809 spike protein sequences were

97    obtained from the "spikeprot0304" file. To determine the number of identical isolate

98    sequences for different spike protein sequences in the GISAID database, the 369809

99    spike protein sequences were further filtered using different spike protein sequences

100   as references. Ultimately, 188 different spike protein sequences, referred to as protein

101   haplotypes, were obtained. Each protein haplotype consisted of at least 100 identical

102   isolate sequences within the set of 369809 spike protein sequences. For each protein

103   haplotype, one representative accession ID (GIS_ISL_) was selected.

104        To obtain the DNA sequences corresponding to the 188 spike protein haplotypes,

105   I downloaded the complete genomes of these haplotypes from GISAID using their

106   accession IDs. The downloaded complete genomes comprised the SARS-CoV-2 DNA

107   sequences. I aligned the 188 complete genomes using MAFFT v.7.450 software [12],

108   using the Wuhan-Hu-1 sequence (GenBank accession number: MN908947.3) as the

109   reference sequence. The resulting alignment contained 189 DNA sequences, including

110   the additional Wuhan-Hu-1 sequence. The spike DNA sequences were cut to a new

111   alignment for phylogenetic and selection analyses.

112      To align the 189 spike DNA sequences, the DNA sequences were first translated

113      into protein sequences using the Bioedit software. The translated protein sequences

114      were then aligned using MAFFT v.7.450 software. Based on the alignment of the

115      protein sequences, the corresponding DNA sequences were aligned using the Dambe

116      software [13].

117      **Reconstruction of phylogenetic tree and median join network**

118      I used the jmodeltest software [14] to determine the best evolutionary model for

119      the alignment of the spike DNA sequences. To reconstruct phylogenetic tree, I

120      conducted maximum likelihood (ML) and Bayesian analyses using IQ-TREE

121      software [15] and MrBayes software [16], respectively. In ML analysis, the statistical

122      support for the tree topology was assessed using 1000 bootstrap replicates. In BA

123      analysis, the parameters of the likelihood model were set as nst = 6 and rate =

124      invgamma, as determined by jmodeltest. The analysis was run for $10^7$ generations,

125      with a sample frequency of 1000 and a burn-in of 2500. The consensus tree with

126      posterior probability was constructed based on 7500 trees.

127      I reconstructed a median-join network based on the 189 spike DNA sequences.

128      The lineages of the spike sequences were assigned according to the Pango-lineage

129      nomenclatures [17] in the GISAID. The median network of the 189 spike DNA

130      haplotypes was constructed with PopART software [18]. To enhance the visualization

131   of different lineages in the phylogenetic tree and median-join network, I used

132   Inkscape and PowerPoint to edit the phylogenetic tree and median-join network. In

133   Inkscape, I assigned different colors to the Pango lineages based on hexadecimal

134   codes, while in PowerPoint, I used the corresponding RGB values to color-code the

135   lineages. These editing steps were performed to facilitate the easy identification and

136   differentiation of the various spike protein lineages in the phylogenetic tree and

137   median-join network.

138       To calculate the genetic distances between the major lineages of SARS-Cov-2,

139   the 189 spike DNA haplotypes were divided into nine major groups: Wuhan-Hu-1,

140   BA.1.1.529.2 (BA.2), BA.1.1.529.4 (BA.4), B.1.1.529.5 (BA.5), XBB.1, XBC, XBF,

141   XBM, and XBZ. The net average distance (the net number of amino acid differences

142   per sequence) was computed for all sequence pairs between these major groups using

143   MEGA11 software [19]. The net average distance between two groups is given by

144       $d_A = d_{XY} - ((d_X + d_Y)/2)$

145   Where, $d_{XY}$ is the average distance between groups X and Y, and $d_X$ and $d_Y$ are the mean

146   within-group distances [19]. The analysis assumed a uniform rate among sites, and

147   pairwise deletion was used to handle gaps between sequences.

148       To determine whether specific amino acid sites in the spike proteins of SARS-

149   Cov-2 were under selection, the nonsynonymous versus synonymous substitution

8

150    ratio (dn/ds ratio = ω) was calculated using the site model in the codeml program of

151    the PAML software [20]. The ω ratio provides information about the balance between

152    nonsynonymous (amino acid-changing) and synonymous (amino acid-preserving)

153    substitutions at each site. A value of ω < 1 suggests purifying (negative) selection, ω =

154    1 suggests neutral evolution, and ω > 1 suggests positive (diversifying) selection.

155    Likelihood ratio tests were performed to compare different evolutionary models:

156    M0 (one ratio) versus M3 (discrete), M1a (nearly neutral) versus M2 (selection), and

157    M7 (beta) versus M8 (beta & ω). The Bayes empirical Bayes method was used to

158    calculate posterior probabilities for site classes [21]. If the likelihood ratio test is

159    statistically significant, it suggests that the amino acid sites are under selection. It is

160    important to note that only the 188 spike DNA haplotypes were analyzed in this study.

161    The Wuhan-Hu-1 sequence was not included in the analyses due to the absence of

162    Wuhan-Hu-1 spike protein haplotypes in the GISAID database from December 1,

163    2012, to February 2013. Amino acid sites with gaps in the spike DNA sequence

164    alignment were deleted because the nonsynonymous versus synonymous substitution

165    value cannot be calculated in the PAML software. The site numbering used the spike

166    protein (protein ID=QHD416.1) of the Wuhan-Hu-1/2019 (GenBank accession

167    number MN908947.3) as the reference for consistency.

168    Results

169 **Characteristics of the spike protein sequences**

170 According to the filtering criteria mentioned, a total of 369809 spike protein

171 sequences were obtained from the spikeprot0304 file. Among these sequences,

172 221323 isolates were collected in December 2022, 119971 isolates in January 2023,

173 and 28515 isolates in February 2023. No isolate was filtered out in March 2023. The

174 number of isolate sequences versus amino acid lengths of spike protein sequences is

175 as follows: 1710 isolate sequences had 1266 amino acids, 57587 isolate sequences

176 had 1267 amino acids, 216386 isolate sequences had 1268 amino acids, 45463 isolate

177 sequences had 1269 amino acids, 47036 isolate sequences had 1270 amino acids, 547

178 isolate sequences had 1271 amino acids, 528 isolate sequences had 1272 amino acids,

179 and 253 isolate sequences had 1273 amino acids. Other spike protein sequences with

180 lengths of 1259, 1260, 1261, 1262, 1263, 1264, 1265, 1274, 1275, 1276, 1277, 1281,

181 1283, or 1319 amino acids had fewer than 72 isolate sequences (Fig. 1). Out of the

182 189 spike protein haplotypes analyzed, there were 4 haplotypes with 1266 amino

183 acids, 36 haplotypes with 1267 amino acids, 106 haplotypes with 1268 amino acids,

184 16 haplotypes with 1269 amino acids, 25 haplotypes with 1270 amino acids, one

185 haplotype with 1272 amino acids, and one haplotype with 1273 amino acids. The

186 haplotype with 1273 amino acids is the Wuhan-Hu-1 sequence, but its spike protein

187 haplotype was not found in the GISAID database from December, 2022 to February,

188    2023.

189    **Net average genetic distances of spike proteins between major lineages of SARS-**

190    **CoV-2**

191      The net average genetic distances of spike protein between Wuhan-Hu-1 and

192    B.1.1.529.2 (BA.2), B.1.1.529.4 (BA.4), B.1.1.529.5 (BA.5), XBB, XBC, XBF, and

193    XBM were 34.54, 31, 37.07, 36.62, 35, 37, 33, and 31 amino acids per sequence,

194    respectively. The net average genetic distances of spike protein between B.1.1.529.2

195    (BA.2) and B.1.1.529.4 (BA.4), B.1.1.529.5 (BA.5), XBB, XBC, XBF, XBM, and

196    XBZ were 9.41, 7.3, 11.87, 16.71, 1.71, 11.41, and 8.67 amino acids per sequence,

197    respectively. The net average genetic distances of spike protein between B.1.1529.4

198    (BA.4) and B.1.1.529.5 (BA.5), XBB, XBC, XBF, XBM, and XBZ were 1.99, 13.18,

199    15, 12, 4, and 4 amino acids per sequence, respectively. The net amino acid

200    differences per sequence of spike protein between B.1.1.529.5 (BA.5) and XBB, XBC

201    XBF, XBM, and XBZ were 11.73, 14.12, 10.52, 3.9, and 1.4 amino acids per

202    sequence, respectively. The net average genetic distances of spike protein between

203    XBB and XBC, XBF, XBM, and XBZ were 19.62, 11.62, 15.18, and 12.93 amino

204    acids per sequence, respectively. The net average genetic distances of spike protein

205    between XBC, and XBF, XBM, and XBZ were 18, 17, 17 amino acids per sequence,

206    respectively. The net average genetic distances of spike protein between XBF and

11

207    XBM and XBZ was 14 and 12 amino acids per sequence, respectively. The net

208    average genetic distances of spike protein between XBM and XBZ was 6 amino acids

209    per sequence (Table.1).

210    **Phylogenetic analyses of spike DNA sequences**

211       The phylogenetic tree of 189 spike DNA sequences (Fig. 2) consisted of three

212    major clades. Clade I consisted of lineages or descendants of BQ.1, BF, and DN. It

213    was positioned closer to the root of the tree. Clade II consisted of lineages or

214    descendants of BA.5. It was located between clade I and clade III in the phylogenetic

215    tree. Clade III was further distal to the root compared to clade II and consisted of

216    subclades A, B, C, and D. Subclade A consisted of lineages or descendants of CM.

217    Subclade B encompasses lineages or descendants of CH.1, CA, CV, and BR.

218    Subclade C consisted of lineages or descendants of BN.1. Subclade D consisted of the

219    lineage or descendant of XBB lineages. In the maximum likelihood (ML) analysis, it

220    was found that the sequences BF.1.1 (EPI_ISL_16152392) and BF.7

221    (EPI_ISL_16080401) within clade I occupied the most basal position when the

222    phylogenetic tree was rooted by the Wuhan-Hu-1 sequence. Statistical analyses,

223    including bootstrap values and posterior probabilities, provided strong support for the

224    monophyly (common ancestry) of clade III and its subclades A, C, and D. A bootstrap

225    value or posterior probability of more than 0.95 indicated a high level of confidence

12

226  in the grouping of sequences within these clades.

**227  Median-join network of spike DNA sequences**

228  Median-join network (Fig. 3) showed that the BF.11 (EPI_ISL_16152392)

229  connected to Wuhan-Hu-1 Spike DNA sequences with 29 nucleotide substitutions.

230  The network can be classified into six major clusters, i.e., BQ.1, BA.5, CH.1.1, CM,

231  BN.1 and XBB.1. The BQ.1 cluster had two diversifying centers. In the BQ.1 cluster's

232  first diversifying center, there were nine haplotypes with the following GISAID

233  accession IDs (EPI_ISL_): 16027638, 16028737, 16029423, 16052382, 16052485,

234  16064186, 16077475, 16113812, and 16660463. It is worth noting that these nine

235  DNA sequences were considered identical in the analysis because the PopART

236  software only counted nucleotide substitutions and did not count insertions or

237  deletions in the alignment. In the BQ.1 cluster's second diversifying center, there were

238  seven sequences with GISAID accession IDs (EPI_ISL_) of 16028751, 16028774,

239  16029345, 16029559, 16052449, 16131848, and 16217334. These sequences also

240  exhibited differences due to insertions and deletions. Among them, the spike protein

241  sequence of EPI_ISL_16028774 was the most abundant, with 16194 isolates recorded

242  in GISAID. The BA.5 cluster consisted of three haplotypes with GISAID accession

243  IDs (EPI_ISL_) of 15973011, 16029234, and 16059569. These three spike DNA

244  sequences also exhibited variations due to insertions and deletions. Among them, the

13

245     spike protein sequence of EPI_ISL_15973011 was the most abundant, with 18098

246     isolates recorded in GISAID. The CH.1.1 cluster consisted of six spike DNA

247     haplotypes that had diversified from an unknown haplotype. Among them, the spike

248     protein haplotype (EPI_ISL_16044651) was the most abundant, with 7329 isolates

249     recorded. It differed from the haplotype of EPI_ISL_16028739 (BA.5.11) by 11

250     nucleotide substitutions. The CM cluster consisted of two haplotypes, namely

251     EPI_ISL_16093062 and EPI_ISL_16029195, with 349 and 857 isolates, respectively.

252     The spike DNA sequence of EPI_ISL_16029195 differed from that of

253     EPI_ISL_16028774 (BQ.1.1) by 13 nucleotide substitutions. The XBB cluster

254     consisted of 14 haplotypes, with its diversifying center consisting of two haplotypes

255     with GISAID accession IDs (EPI_ISL_) of 16044705 and 16206019. Among these

256     haplotypes, the spike protein haplotype of EPI_ISL_16044705 was the most

257     abundant, with 24144 isolates recorded. It differed from the haplotype of

258     EPI_ISL_16040256 (BN.1.2) by 13 nucleotide substitutions and from the haplotype

259     of EPI_ISL_16028739 (BA.5.11) by 22 substitutions. The DNA haplotype of

260     EPI_ISL_16168343 (XBC.1) differed from that of EPI_ISL_15973011 (BA.5.2) by

261     19 nucleotide substitutions.

262     **Positive selection sites of spike protein**

263    The values of likelihood ratio tests of M0 versus M3, M1a versus M2, and M7

264    versus M8 comparisons were larger the critical values at 0.01 level. The results

265    suggest that the M3, M2, and M8 models were statistically better than M0, M1a, and

266    M7 models, respectively. The Bayes empirical Bayes (BEB) analyses of M2 models

267    identified the 25 amino-acid sites under positive selection. These sites were located at

268    the positions of L5**, W152 **, G181 **, N185 **, G213*, V213*, H245*, Y248*,

269    D253**, S255*, S256**, G257*, R346**, R408*, K444**, V445**, G446**,

270    N450**, L452**, N460*, F486**, Q613**, Q675*, T883**, P1162**, and V1264**,

271    which were statistically significant at 0.05 (*) and 0.01 (**) levels. The M8 model

272    identified an additional one more site at V83* which was not identified by M2 model

273    (Table 2). The site of L5 was located in signal peptide domain (SP) of the spike

274    protein. The V83, W152, G181, N185, G213, H245, Y248, D253, S255, S256 and

275    G257 were located in N-terminal domain (NTD). The R346, R408, K444, V445,

276    G446, N450, L452, N460, and F486 were located in receptor binding domain (RBD).

277    The Q613 and Q675 were located in C-terminal domain 2 (CTD2). The T883 was

278    located in fusion-peptide proximal region (FPPR). The P1162 was located between

279    HR1 and HR2. The V1264 was located in cytoplasmic tail (CT). The nonsynonymous

280    substitutions of selection sites ranged from 4 to 11 in each protein haplotype and the

281    same nonsynonymous substitution in the same selection sites usually occurred in

15

282    different lineages except the substitutions of V83A, V213E, and V445P were

283    exclusively occurred in the XBB lineage. Among these selection sites, the site of 444

284    had the largest amino acid diversity. The nonsynonymous substitutions included

285    K444R, K444T, K444M, and K444N that were occurred in 7, 94, 4, 5 of 188 protein

286    haplotypes, respectively.

287    **Discussion**

288    The presence of long and short spike protein sequences, with variations in amino

289    acid length compared to the original Wuhan-Hu-1 spike protein, is not completely due

290    to incomplete sequencing or sequence error. This conclusion is based on several

291    observations made in the study. Firstly, the sequences did not contain ambiguous

292    amino acids (represented by X). Secondly, all the sequences analyzed contained a

293    start codon (M), Additionally, most sequences had a complete C-terminus domain. It

294    is important to note that these variations in amino acid length were typically observed

295    in the signal peptide or N-terminus domains, and rarely in the S2 region. Importantly,

296    these insertions and deletions were never observed in the receptor binding domain

297    (RBD) of the spike protein. The RBD is responsible for binding to the ACE-2

298    receptor, which is essential for viral entry into host cells. The fact that strains with

299    long or short spike proteins still maintained infectivity suggests that they were still

300    able to bind to the ACE-2 receptor despite these variations of sequence lengths.

16

301    The results showed that the net average genetic distances of spike protein between

302    the Wuhan-Hu-1 strain and lineages of B.1.1.529.2, B.1.1.529.4, B.1.1.529.5, XBB,

303    XBC, XBF, XBM, and XBZ ranged from 30.07 (between Wuhan-Hu-1 and

304    BA.1.1.529.5) to 37 (between Wuhan-Hu-1 and XBF) amino acids per sequence. The

305    results showed there was a great difference between the original (Wuhan-Hu-1) and

306    current strains. Furthermore, the study specifically mentions the genetic distances

307    between the XBB strain and several other strains. The genetic distances between XBB

308    and lineages of B.1.1.529.2, B.1.1.529.4, B.1.1.529.5, XBC, XBF, XBM, and XBZ

309    were 11.87, 13.18, 11.73, 19.62, 11.62, 15.18, and 12.93, respectively. Among these

310    strains, XBC had the largest difference from XBB, with 19.62 amino acids per

311    sequences. XBC was a recombinant of BA.2 Omicron (the most mutated) and

312    B.1.617.2 Delta (the most severity) strains [22, 23]. It is important to continue

313    surveillance and monitor the evolution of XBC.

314    The results of phylogenetic tree (Fig. 2) and median-join network (Fig. 3) revealed

315    that the presence of multiple lineages of SARS-CoV-2 during December 2022 to

316    February 2023. However, the majority of these lineages were descendants of three

317    major lineages: BA.2, BA.5, and XBB. To help summarize the relationships between

318    the lineages, the study employed the use of simplified names based on the Pango

319    lineage nomenclature. However, the full names providing a more detailed and precise

320    identification of the lineages.

321    Firstly, the BA.2 (BA.1.1.529.2) consisted of the sub-lineages of CM

322    (B.1.1.529.2.3.20), CA (B.1.1.529.2.75.2), CV (B.1.1.529.2.75.3.1.1.3), DV

323    (B.1.1.529.2.75.3.4.1.1.1.1.1), CH (B.1.1.529.2.75.3.4.1.1), BR (B.1.1.529.2.75.4),

324    BN (B.1.1.529.2.75.5), EJ.2 (B.1.1.529.2.75.5.1.3.8.2) and BY (B.1.1.529.2.75.6).

325    Secondly, the BA.5 (BA. 1.1.529.5) consisted of eight major sub-lineages, i.e.,

326    BA.5.1, BA. 5.2, BA.5.3, BA.5.5, BA.5.6, BA.5.9, BA.5.10, and BA.5.11 in this

327    study. The descendants of BA.5.1 (B.1.1.529.5.1) consisted of BA.5.1.5

328    (B.1.1.529.5.1.5), BA.5.1.12 (B.1.1.529.5.1.12), BA.5.1.27 (B.1.1.529.5.1.27), and

329    CL.1 (B.1.1.529.5.1.29.1) in this study. The descendants of BA.5.2 consisted of

330    BA.5.2.1 (B.1.1.529.5.2.1), BF.5 (B.1.1.529.5.2.1.5), BF.7 (B.1.1.529.5.2.1.7), BU.1

331    (B.1.1.529.5.2.16.1), CR.1.1 (B.1.1.529.5.2.18.1.1), CR.1.2 (B.1.1.529.5.2.18.1.2),

332    CN.1 (B.1.1.529.5.2.21.1), CN.2 (B.1.1.529.5.2.21.2), BA.5.2.23 (B.1.1.529.5.2.23),

333    CK.2 (B.1.1.529.5.2.24), in this study.   The descendants of BA.5.3 (B.1.1.529.5.3)

334    consisted of BQ.1 (B.1.1.529.5.3.1.1.1.1.1), DU.1 (B.1.1.529.5.3.1.1.1.1.1.1.2.1), and

335    CQ (B.1.1.529.5.3.1.4.1.1) in this study. The descendants of BA.5.6 (B.1.1.529.5.6.)

336    consisted of BW.1.1 (B.1.1.529.5.6.2.1.1) in this study. The descendants of BA.5.10

337    ((B.1.1.529.5.10) consisted of DF (B.1.1.529.5.10.1) in this study. The BA.5.11

338   consisted of the BA.5.11 only. Thirdly, XBB was the recombinant of two BA.2

339   lineages, i.e., BJ.1 and BM1.1.1 [24]. The EG.1 and FL.10 were the abbreviations of

340   XBB.1.9.2.1 and XBB.1.9.1.10, respectively. The other recombinants include XBC (a

341   recombinant of BA.2 Omicron and Delta), XBF (a recombinant of BA.5 and

342   BA.2.75), and XBZ (a recombinant of BA.5.2 and EF.1.3) based on the Covid-lineage

343   Pango designation (Roemer, 2022) [23]. The most dominant variant was the strain

344   BQ. 1.1.23 with the representative accession number of EPI_ISL_16027638, and had

345   55919 identical isolate sequences, following by XBB.1.5 (representative accession

346   number EPI_ISL_16044705, 24133 identical isolate sequences), and BA.5.11

347   (representative accession number EPI_ISL_16028739, 21798 identical isolate

348   sequences) during December, 2022 to February, 2023.

349     The previous study demonstrated that certain mutations in the receptor-binding

350   domain (RBD) of the spike protein, specifically at positions R346, K356, K444,

351   V445, G446, N450, L452, N460, F486, F490, R493, or S494, could lead to the

352   evasion of neutralizing monoclonal antibodies (mAbs) or enhance binding to the

353   ACE2 receptor (Cao et al., 2022). In the present study, we found that mutations at

354   R346, K444, V445, G446, N450, L452, N460, and F486 had a nonsynonymous

355   versus synonymous substitution ratio greater than 1, indicating positive selection. This

356   suggests that these sites were undergoing evolutionary changes that may confer

357     selective advantages to the virus. However, mutations at K356, F490, F493, or S494

358     did not exhibit a nonsynonymous versus synonymous substitution ratio greater than 1

359     in the present analysis, suggesting that these sites were not under positive selection

360     during the specific time frame examined (December 2022 to February 2023) in the

361     study (Table 2). This finding contrasts with the previous study, which analyzed

362     sequences from January 2021 to October 2022. I propose that the discrepancy in

363     results between the previous and present studies may be attributed to antigenic shift.

364     It's possible that the evolutionary dynamics and selective pressures acting on SARS-

365     CoV-2 may have shifted, leading to different mutations being favored in different

366     time periods. Additionally, the present study identified positive selection for

367     mutations occurring outside of the RBD domain. These sites included L5, V83,

368     W152, G181, N185, G213, H245, Y248, D253, S255, S256, Q613, Q675, T883,

369     P1162, and V1264. However, the effects of these mutations on the fitness of SARS-

370     CoV-2 remain to be investigated.

371         In this study, it was observed that multiple strains coexisted between December

372     2022 and February 2023. However, the majority of these strains belonged to the

373     lineages or sub-lineages of BA.2 (BA.1.1.529.2), BA.5 (BA.1.1.529.5), and XBB

374     (Fig. 2). The diversifying centers of BN.1.2, BQ.1, BA.5.11, XBB were the isolate

375     sequences with representative accession IDs (EPI_ISL_) of 16040256, 16027638,

376    16028739, and16044705, respectively (Fig. 3). I propose that the complete sequences

377    or the receptor binding domain of these spike DNA sequences could be potential

378    candidates for vaccine design. This suggests that these sequences may possess

379    important characteristics that can be utilized in the development of effective vaccines

380    against SARS-CoV-2.

381        As of June 10, 2023, just before submitting our manuscript, the XBB.1.5 and

382    XBB.1.16 strains have emerged as the globally dominant strains, with respective

383    frequencies of 72% and 12% based on data from GISAID [25]. These strains have

384    gained prominence and become widespread within the population. Additionally, the

385    XBC variant is a recombinant of BA.2 (Omicron) and B.1.617.2 (Delta) [17, 23].

386    XBC exhibits significant differences from the XBB lineages and its sub-lineages,

387    making it a distinct variant from XBB. Considering the success of the Omicron

388    variant [3, 4], I propose that the XBC.1 strain or its sub-lineages could potentially

389    become dominant strains following the XBB.1 lineage and its sub-lineages. Continued

390    surveillance and research are necessary to monitor the evolution and potential impact

391    of these variants on public health.

392

393 **Acknowledgements**

394     We gratefully acknowledge all data contributors, i.e., the Authors and their

395 Originating laboratories responsible for obtaining the specimens, and their

396 Submitting laboratories for generating the genetic sequence and metadata and

397 sharing via the GISAID Initiative, on which this research is based (Supplementary

398 Table 1).

399

400 **Funding**

401 Not applicable

402

403 **Availability of data and materials**

404 All sequences were downloaded from Global Initiative on Sharing

405 Avian Influenza Data (GISAID, https://www.gisaid.org/) and GenBank (https://

406 www.ncbi.nlm.nih.gov/nucleotide/).

407

408 **Ethics approval and consent to participate**

409 Not applicable.

410

411 **Consent for publication**

412    Not applicable.

413

414    **Competing interests**

415    The authors declare that they have no competing interests.

416

417    **Author details**

418    [1]Department of Life Sciences, National Chung Hsing University, Taiwan, R.O.C.

419    145 Xingda Road., South District., Taichung City 40227, Fax: +886-4-22874740

420    Tel: Tel：+886-4-22840416

421

422    *Correspondence: Hsiao-Wei Kao

423    E-mail: hkao@dragon.nchu.edu.tw

424

425     **Figure legends**

426

427     **Figure 1**. Number of isolate sequences in GISAID of different lengths of spike

428     protein from December 2022 to February 2023.

429

430     **Figure. 2**. Phylogeny of SARS-CoV-2 spike DNA sequences. The terminal node

431     (leaf) is the GISAID ID of the sequence followed by the lineage name in

432     parentheses, the length of the spike protein, and the number of isolates. Statistical

433     supports are labeled on the branches. The values below 60% are not labeled.

434

435     **Figure. 3.** Median-join network of SARS-CoV-2 spike DNA sequences from

436     December 2022 to February 2023. GISAID ID was labeled inside the circles. The

437     number of isolates and lineages were labeled outside the circles. The number of

438     nucleotide substitutions between haplotypes was labeled on the lines with hatch

439     bars. When the hatch bars exceed 5, the substitutions were also labeled with

440     numbers.

441

442

24

## References

443    1. World Health Organization. WHO Statement on the fifteenth meeting of the IHR

445      (2005) Emergency Committee on the COVID-19 pandemic. World Health

446      Organization, Geneva, Switzerland. 2023. https://www.who.int/news/item/05-05-

447      2023-statement-on-the-fifteenth-meeting-of-the-international-health-regulations-

448      (2005)-emergency-committee-regarding-the-coronavirus-disease-(covid-19)-

449      pandemic

450    2. World Health Organization. Update on Omicron. World Health Organization,

451      Geneva, Switzerland. 2021. https://www.who.int/news/item/28-11-2021-update-

452      on-omicron

453    3. Martin DP, Lytras S, Lucasi AG, Maier W, Grüning B, Shank SD et al. Selection

454      analysis identifies cluster of unusual mutational changes in omicron lineage BA.1

455      that likely impact spike function. Mol Biol Evol. 2022;39 (4): msac061.

456      https://doi.org/10.1093/molbev/msac061

457    4. Dejnirattisai W, Huo J, Zhou D, Zahradník J, Supasa P, Liu C et al. SARS-CoV-2

458      Omicron-B.1.1.529 leads to widespread escape from neutralizing antibody

459      responses. Cell 2022;185(3): 467-484. https://doi.org/10.1016/j.cell.2021.12.046

460    5. Tegally H, Moir M, Everatt J, Giovanetti M, Scheepers C, Wilkinson E et al.

461      Emergence of SARS-CoV-2 Omicron lineages BA.4 and BA.5 in South Africa.

25

462    Nat Med. 2022;28:1785–1790. https://doi.org/10.1038/s41591-022-01911-2

463    6.  Willett BJ, Grove J, MacLean OA, Wilkie C, Lorenzo GD, Furnon W, et al.

464    SARS-CoV-2 Omicron is an immune escape variant with altered cell entry

465    pathway. Nat. Microbiol. 2022;7:1161-1179. https://doi.org/10.1038/s41564-022-

466    01143-7

467    7. Wang L, Møhlenberg M, Wang P, Zhou H. Immune evasion of neutralizing

468    antibodies by SARS-CoV-2 Omicron. Cytokine Growth Factor Rev. 2023;70:13–25.

469    8. Cao Y, Jian F, Wang J, Yu Y, Song W, Yisimayi A, et al. Imprinted SARS-CoV-2

470    humoral immunity induces convergent Omicron RBD evolution. Nature

471    2023;614:521–529. https://doi.org/10.1038/s41586-022-05644-7

472    9. Yang, Z., Nielsen, R., Goldman, N. and Perdersen, A.M. Condon-substitution

473    models for heterogeneous selection pressure at amino acid sites. Genetics

474    2000;155(1): 431-449. https://doi.org/10.1093/genetics/155.1.431

475    10. Elbe S, Buckland-Merrett G. Data, disease and diplomacy: GISAID's

476    innovative contribution to global health. Global Chall. 2017;1(1):33-46.

477    https://doi.org/10.1002/gch2.1018

478    11. Hall TA. BioEdit: A user-friendly biological sequence alignment editor and

479    analysis program for Windows 95/98/NT. Nucleic Acids Symp. Ser 1999,41:95-98.

480    12. Katoh K, Standley DM. (2013). MAFFT multiple sequence alignment software

481    version 7: Improvements in performance and usability. Mol. Biol. Evol.

482    2013;30(4):772-780. https://doi.org/10.1093/molbev/mst010

483    13. Xia X. DAMBE7: New and improved tools for data analysis in molecular biology

484    and evolution. Mol. Biol. Evol. 2018;35(6):1550–1552.

485    https://doi.org/10.1093/molbev/msy073

486    14. Posada D. Jmodeltest: Phylogenetic model averaging. Mol Biol Evol.

487    2008;25(7):1253-1256. https://doi.org/10.1093/molbev/msn083

488    15. Nguyen LT, Schmidt HA, Haeseler A, Minh BQ. IQ-TREE: A fast and effective

489    stochastic algorithm for estimating maximum likelihood phylogenies. Mol. Biol.

490    Evol. 2015;32(1):268-274. https://doi.org/10.1093/molbev/msu300

491    16. Huelsenbeck, JP, Ronquist F. MRBAYES: Bayesian inference of phylogenetic

492    tree. Bioinformatics 2001;17(8):754-755.

493    https://doi.org/10.1093/bioinformatics/17.8.754

494    17. Rambaut A, Holmes EC, O'Toole Á, Hill V, McCrone JT, Ruis C et al. A

495    dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic

496    epidemiology. Nat Microbiol. 2020;5,1403–1407. https://doi.org/10.1038/s41564-

497    020-0770-5

498    18. Leigh JW, Bryant D. POPART: full-feature software for haplotype network

499    construction. Methods Ecol Evol. 2015;6 (9):1110–1116.

500    https://doi.org/10.1111/2041-210X.12410

501    19. Tamura K, Stecher G, Kumar S. MEGA 11: Molecular evolutionary genetics

502    analysis Version 11. Mol Biol Evol. 2021;38(7):3022-3027.

503    https://doi.org/10.1093/molbev/msab120

504    20. Yang Z. PAML 4: Phylogenetic analysis by maximum likelihood.

505    Mol. Biol. Evol. 2007;24(8):1586–1591.

506    https://doi.org/10.1093/molbev/msm088

507    21. Yang, Z., Wong, W. S., and Nielsen, R. (2005). Bayes empirical Bayes inference

508    of amino acid sites under positive selection. Mol. Biol. Evol. 2005;22(4):1107–

509    1118. https://doi.org/10.1093/molbev/msi097

510    22. Varea-Jiménez E, Cano EA, Vega-Piris L, Sánchez EVM, Mazagatos C,

511    Rodríguez-Alarcón LGSM et al. Comparative severity of COVID-19 cases

512    caused by Alpha, Delta or Omicron SARS-CoV-2 variants and its association

513    with vaccination, Enferm Infecc Microbiol Clin.2022;

514    https://doi.org/10.1016/j.eimc.2022.11.003

515    23. Roemer C. Cov-Lineages/Pango-designations-Lineage description. Github: cov-

516    lineages/pango-designation. 2022. https://github.com/cov-lineages/pango-

517    designation

518    24. Tamura T, Ito J, Uriu K, Zahradnik J, Kida I, Anraku Y et al. Virological

519      characteristics of the SARS-CoV-2 XBB variant derived from recombination of

520      two Omicron subvariants. Nat Commun. 2023;14:2800.

521      https://doi.org/10.1038/s41467-023-38435-3

522   25. Hadfield J, Megill C, Bell, SM, Huddleston J, Potter B, Callender C. et al. (2018)

523      Nextstrain: real-time tracking of pathogen evolution. *Bioinformatics*

524      2018;34(23):4121-4123. https://doi.org/10.1093/bioinformatics/bty407

525

526  **Table 1**. The net average genetic distances of per sequence between nine-lineage

527  spike proteins. All ambiguous positions were removed for each sequence pair

528  (pairwise deletion option).

|  | Wuhan-Hu-1 | B.1.1.529.2 (BA.2) | B.1.1.529.4 (BA.4) | B.1.1.529.5 (BA.5) | XBB | XBC | XBF | XBM |
|---|---|---|---|---|---|---|---|---|
| B.1.1.529.2 (BA.2) | 34.54 | | | | | | | |
| B.1.1.529.4 (BA.4) | 31.00 | 9.41 | | | | | | |
| B.1.1.529.5 (BA.5) | 37.07 | 7.30 | 1.99 | | | | | |
| XBB | 36.62 | 11.87 | 13.18 | 11.73 | | | | |
| XBC | 35.00 | 16.71 | 15.00 | 14.12 | 19.62 | | | |
| XBF | 37.00 | 1.71 | 12.00 | 10.52 | 11.62 | 18.00 | | |
| XBM | 33.00 | 11.41 | 4.00 | 3.90 | 15.18 | 17.00 | 14.00 | |
| XBZ | 31.00 | 8.67 | 4.00 | 1.40 | 12.93 | 17.00 | 12.00 | 6.00 |

529

530 **Table 2**. Likelihood ration test of M0 vs M3, M1a vs M2, M7 vs M8, and amino acid
531 site of spike protein under positive selection.

| Parameter | M0 | M3 | M1a | M2 | M7 | M8 |
|---|---|---|---|---|---|---|
| -lnL | 7508.42 | 7281.49 | 7414.25 | 7295.08 | 7421.13 | 7322.99 |
| 2ln (L1-L0) | 453.86 (between M0 and M3) | | 238.34 (between M1a and M2) | | 196.28 (between M7 and M8) | |
| df between models | 4 | | 2 | | 2 | |
| Chi square test | P<0.01 | | P<0.01 | | P<0.01 | |
| Positive selective sites | Not allow | Not allow | Not allow | L5**, W152**, G181**, N185*, V213*, H245*, Y248*, D253**, S255*, S256**, G257**, R346**, R408**, K444**, V445**, G446**, N450**, L452**, N460*, F486**, Q613**, Q675*, T883**, P1162**, V1264** | Not allow | L5**, V83*, W152 **, G181 **, N185 *, V213*, H245*, Y248*, D253**, S255*, S256**, G257*, R346**, R408*, K444**, V445**, G446**, N450**, L452**, N460*, F486**, Q613**, Q675*, T883**, P1162**, V1264** |

532 * Statistically significant at 0.05, ** statistically significant at 0.01.
533

534



535

**Fig 1**. Number of isolate sequences versus different lengths of spike protein in GISAID from December 2022 to February 2023.

538

539



**Fig. 2.** Phylogeny of SARS-CoV-2 spike DNA sequences. The terminal node (leaf) is the GISAID ID of the sequence followed by the lineage name in parentheses, the length of the spike protein, and the number of isolates. Statistical supports are labeled on the branches. The values below 60% are not labeled.

**Fig. 3.** Median-join network of SARS-CoV-2 spike DNA sequences from December 2022 to February 2023. GISAID ID was labeled inside the circles. The number of isolates and lineages were labeled outside the circles. The number of nucleotide substitutions between haplotypes was labeled on the lines with hatch bars. When the hatch bars exceed 5, the substitutions were also labeled with numbers.