1    **Genomic surveillance unfolds the dynamics of SARS-CoV-2 transmission and**

2    **divergence in Bangladesh over the past two years**

3

4    Tushar Ahmed Shishir [a, b], Taslimun Jannat[a], Iftekhar Bin Naser[a]#

5

6    [a]Department of Mathematics and Natural Sciences, BRAC University, Dhaka – 1212,

7    Bangladesh.

8    [b]Rangamati General Hospital, Rangamati – 4500, Chattogram, Bangladesh.

9

10   **Running Head:** SARS-CoV-2 transmission and divergence in Bangladesh

11

12   #Address correspondence to:

13   Iftekhar Bin Naser,

14   Email: iftekhar.naser@bracu.ac.bd

15

16   **Word count –**

17   1.  Abstract – 187

18   2.  Importance – 147

19   3.  Introduction – 738

20   4.  Results – 2245

21   5.  Discussion – 874

22   6.  Materials and Methods – 346

1

## Abstract

The highly pathogenic virus SARS-CoV-2 has shattered the healthcare system of the world causing the COVID-19 pandemic since first detected in Wuhan, China. Therefore, scrutinizing the genome structure and tracing the transmission of the virus has gained enormous interest in designing appropriate intervention strategies to control the pandemic. In this report, we examined 4622 sequences from Bangladesh and found that they belonged to thirty-five major PANGO lineages, while Delta alone accounted for 39%, and 78% were from just four primary lineages. Our research has also shown Dhaka to be the hub of viral transmission and observed the virus spreading back and forth across the country at different times by building a transmission network. The analysis resulted in 7659 unique mutations, with an average of 24.61 missense mutations per sequence. Moreover, our analysis of genetic diversity and mutation patterns revealed that eight genes were under negative selection pressure to purify deleterious mutations, while three genes were under positive selection pressure.

## Importance

With 29,122 deaths, 1.95 million infections and a shattered healthcare system from SARS-CoV-2 in Bangladesh, the only way to avoid further complications is to break the transmission network of the virus. Therefore, it is vital to shedding light on the transmission, divergence, mutations, and emergence of new variants using genomic data analyses and surveillance. Here, we present the geographic and temporal distribution of different SARS-CoV-2 variants throughout Bangladesh over the past two years, and their current prevalence. Further, we have developed a transmission network of viral spreads, which in turn will help take intervention measures. Then we analyzed all the mutations that occurred and their effect on evolution as well as the currently present mutations that could trigger a new variant of concern. In short, together with an ongoing genomic surveillance

46  program, these data will help to better understand SARS-CoV-2, its evolution, and pandemic

47  characteristics in Bangladesh.

48

49  **Keywords:** SARS-CoV-2; COVID-19; Genetic diversity; Molecular surveillance; Evolution;

50  Pandemic

51

52  **Introduction**

53  The Coronavirus disease (COVID-19) pandemic was initially reported as an unknown respiratory

54  illness towards the end of 2019. However, it eventually became evident that a novel SARS-like

55  coronavirus was causing the infections and the virus was termed severe acute respiratory syndrome

56  coronavirus 2 (SARS-CoV-2) (1). Originating in Wuhan, China, SARS-CoV-2 has spread across

57  220 countries and territories, infecting 488.75 million and causing the death of 6.17 million people

58  till 31st March 2022, resulting in a global economic crisis, which is the third zoonotic virus after

59  MERS-CoV and SARS-CoV in 2012 and 2002 respectively (2, 3). The novel virus belonging to

60  the Betacoronavirus genus and Coronaviridae family is a positive-sense, single-stranded ~30 kb

61  long RNA virus. Its genome contains 38% GC content (4), prefers pyrimidine-rich codons over

62  purines (5) and is organized into 11 open reading frames expressing 12 proteins, including two

63  polypeptides, four structural proteins and other accessory proteins (6). Phylogenetically, the virus

64  shares 96% identity with the strain BatCoV RaTG13 of Rhinolophus affinis, and genome

65  sequences along with epidemiological data suggest that SARS-CoV-2 is primarily transmitted

66  from bats to humans (3, 4, 7). A complete genome sequence of the virus was deposited in GenBank

67  on 5th January (NC_045512.2) (8), followed by the submission of 9.74 million complete

68  sequences to GISAID by 25th March 2022 (9).

3

69     According to recent data from Worldometer, the most infected countries are the USA, India, and

70     Brazil, with more than 29, 43, and 33 million cases of infection and thousands of deaths (2, 10).

71     Since the first case was confirmed in Bangladesh on 8th March 2020, there have been 1.95 million

72     positives and 29,122 deaths reported until 31st March 2022 (11). Having such a large population

73     makes Bangladesh more vulnerable to viral transmission, and it is labelled as the second-most

74     infected nation in the South Asian region (10), despite the government imposing lockdowns, social

75     distancing rules and mask mandates to control the situation. Therefore, it is crucial to shed light

76     on the transmission and evolution of the virus inside the country to reduce the fatality where

77     genomic data analyses and surveillance comes into play, which can deliver immense information.

78     Child Health Research Foundation published the first SARS-CoV-2 genome sequence from

79     Bangladesh on 12th May 2020 (12), followed by 5146 further sequences until 25th March 2022

80     (9).

81     To date, Bangladesh has been affected by three waves of COVID-19 with variants of concerns

82     (VOC), including Alpha, Beta, Delta, and Omicron (9). VOC is the name given to a variant of the

83     SARS-CoV-2 virus that has mutations in the spike protein receptor-binding domain, increasing

84     binding affinity within the RBD-hACE2 complex and increasing viral transmission (13, 14).

85     Consequently, the mutations are essential for studying since they alter the antigenic potentials of

86     the epitopes and consequently affect pathogenicity, infectivity, transmissibility, and evading host

87     immunity. SARS-CoV-2 encodes an exoribonuclease that proofreads the errors during viral RNA

88     synthesis; therefore, it has a lower mutation rate than other RNA viruses, which aids in enhancing

89     its ability to adapt to their environment (15, 16). Nevertheless, the virus is accumulating mutations

90     across its genome, leading to the emergence of different variants over time. These mutations are

91     not evenly distributed; for example, some genes are more prone to mutations than others are, and

4

92  cytosine to uracil substitution is more common in SARS-CoV-2, reforming the

93  transition/transversion ratio, which is negatively correlated with evolutionary time (17).

94  Additionally, a variable vaccination rate among the countries increases the risk of SARS-CoV-2

95  mutating into a strain that is resistant to current vaccines and therapies. Consequently, it is essential

96  to continue investigating the mutations of SARS-CoV-2 in order to develop further effective

97  vaccines and therapies, improve pandemic response, and reduce the impact of the pandemic on

98  healthcare and clinical processes.

99  To the best of our knowledge, most previous studies in Bangladesh addressed lineages distribution,

100  source determination, and potential mutations with only a few sequences from the early phase of

101  the outbreak (18, 19). Therefore, in this work, we comprehensively analyzed 4622 whole-genome

102  sequences of SARS-CoV-2 isolated from Bangladesh until 25th March 2022 to understand the

103  distribution of variants and mutation accumulation trends over the year. We have thoroughly

104  studied the temporal and geographical distribution of different lineages inside Bangladesh and

105  built the transmission network to trace their back and forth circulation. To better understand the

106  evolutionary dynamics of SARS-CoV-2 in Bangladesh over the last two years, we examined the

107  genetic diversity among strains, gene-wise mutation distribution, and selection pressures.

108

109  **Results**

110  **SARS-CoV-2 lineage dynamics**

111  To understand the diversity and transmission of the virus, we have confined and analyzed

112  sequences from all administrative divisions of Bangladesh. There were 5146 sequences submitted

113  in GISAID till 25th March 2022, but many of them were incomplete and lacked quality. Therefore,

114  we filtered the sequences based on their completeness, coverage, and gaps, resulting in 4892

115    sequences for downstream analysis (Supplementary file 1). However, when we examined PANGO

116    Lineages, we observed that there were 93 lineages, many of which carried extremely low numbers

117    of the sequences. Hence, we further filtered the sequences and kept only the lineages containing

118    at least ten sequences, resulting in 4622 sequences from 35 lineages. Overall, in the beginning, the

119    country had strains that belonged to the fewest number of PANGO lineages, but this has changed

120    over time (Supplementary file 2). Selected sequences belonging to thirty-five different PANGO

121    lineages provided us with invaluable insight regarding patterns of pandemic and viral spread

122    (Supplementary file 2). As an example, 78% of the sequences were grouped into four lineages,

123    where Delta (B.1.617.2) and its three major sub-lineages (AY.X) combined made up the highest

124    39% of the total sequences, while 20 out of thirty-five lineages held only 9% of sequences even

125    after we filtered out the lineages with very few sequences. The top ten most prevalent lineages

126    were found to be B.1.617.2 (23.34%), B.1.1.25 (20.68%), BA.2 (8.57%), B.1.351.3 (7.27%), B.1.1

127    (2.40%), B.1.1.7 (1.86%), B.1 (1.80%), B.1.351 (1.02%), B.1.36.16 (0.93%) and B.1.1.318
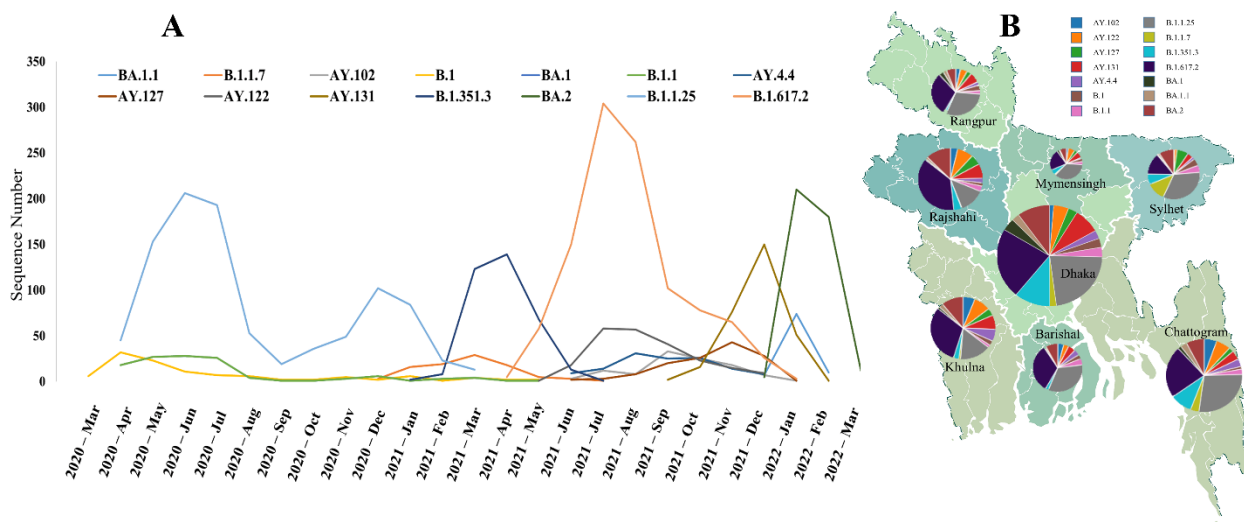
128    (0.74%).

129

130    **Temporal distribution of major lineages**

131    We found that Bangladesh was afflicted by a large number of viruses from 35 different lineages,

132    with the highest diversification occurring between July and September of 2021 with sequences

133    from 20 to 23 lineages (supplementary file 2). The early phase of the pandemic in Bangladesh was

134    started by the introduction of lineage B.1 in March 2020. Multiple occurrences of the introduction

135    of COVID-19 from different countries have previously been reported; for instance, Dhaka was

136    first exposed to COVID-19 with strains from the United Kingdom, while Chattogram was exposed

137    to strains from Saudi Arabia (18). The early phase of the pandemic was generally dominated by

138   imported strains from outside countries, but as the pandemic progressed, mutations changed

139   dynamics and the linage B.1.1.25 took over, with B.1 gradually declining (Fig 1A). B.1.1.25 was

140   the highest prevalent strain until January 2021. Later, the Beta variant (B.1.351) was reported in

141   November 2020, followed by the Alpha variant (B.1.1.7) in December 2020. The B.1.1.7 lineage

142   started talking over the B.1.1.25 lineage following its introduction. This linage was the most

143   frequently detected variant in February 2021, while Beta variants were very less numerous. Despite

144   this, a sub-lineage of Beta variants (B.1.351.3) emerged and outnumbered the Alpha variant in

145   March 2021 (Supplementary file 2). However, the dominance of B.1.351.3 did not last long due

146   to the introduction of the deadly delta variant (B.1.617.2).

147



148

149   **Fig 1: Distribution of major lineages in Bangladesh. A.** Geographic distribution of major

150   lineages at eight administrative divisions of Bangladesh. Dhaka contained the highest number of

151   sequences and maximum diversity, while Mymensingh was the least diverse zone. **B.** Temporal

152   distribution of major lineages. Maximum diversity was observed after the introduction of the Delta

153   variant due to the emergence of different sub-lineage of it. Three major peaks depict the three

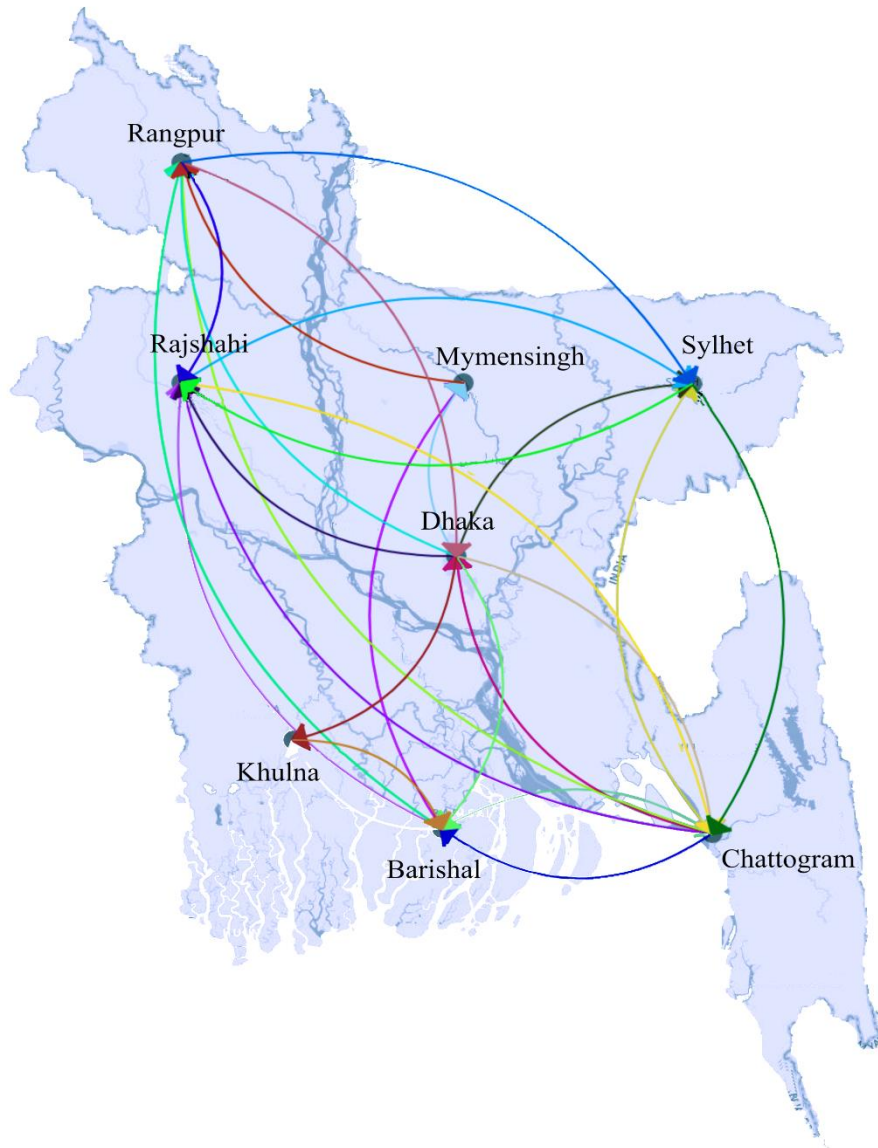154   variants responsible for three COVID-19 waves in Bangladesh.

7

155

156    According to our analysis, B.1.617.2 was the most dominant strain within a month after its

157    introduction in April 2021. A number of distinct A lineages have also been observed, which were

158    mostly sub-lineages of the delta variant, possibly due to the increased transmissibility of the

159    variant. Specifically, AY.122 increased significantly from September 2021 while B.1.617.2 was

160    declining. Meanwhile, the AY.131 lineage first appeared in Bangladesh in October 2021 and

161    dominated all other variants in November 2021; more than half the sequences of December 2021

162    came from this lineage. This variant was eventually replaced by another highly transmissible

163    variant called Omicron (B.1.1.529). Omicron first emerged in Bangladesh in December 2021 and

164    took over within a month. Throughout Bangladesh, the Omicron variant has been dominant since

165    January 2022.

166

167    **Regional distribution of different lineages.**

168    We then conducted a chronological lineages dynamics analysis in order to determine whether the

169    variants were distributed evenly across Bangladesh's administrative divisions. In terms of

170    geographical distribution, Dhaka had the most diversified sequences from all thirty-five lineages,

171    followed by Chattogram from 31. On the contrary, Mymensingh and Rangpur were less diverse

172    areas with sequences from only 22 and 24 lineages, respectively, where most of the lineages

173    represented only one or two sequences (Fig 1B). The Alpha variant was first detected in the Sylhet

174    division and has since spread to the other five divisions except for Barishal and Rangpur, where

175    the Delta and Omicron variant first appeared in Dhaka. Overall, the ratio of the dominating

176    lineages was similar throughout the country, and our analyzed transmission network reflects that
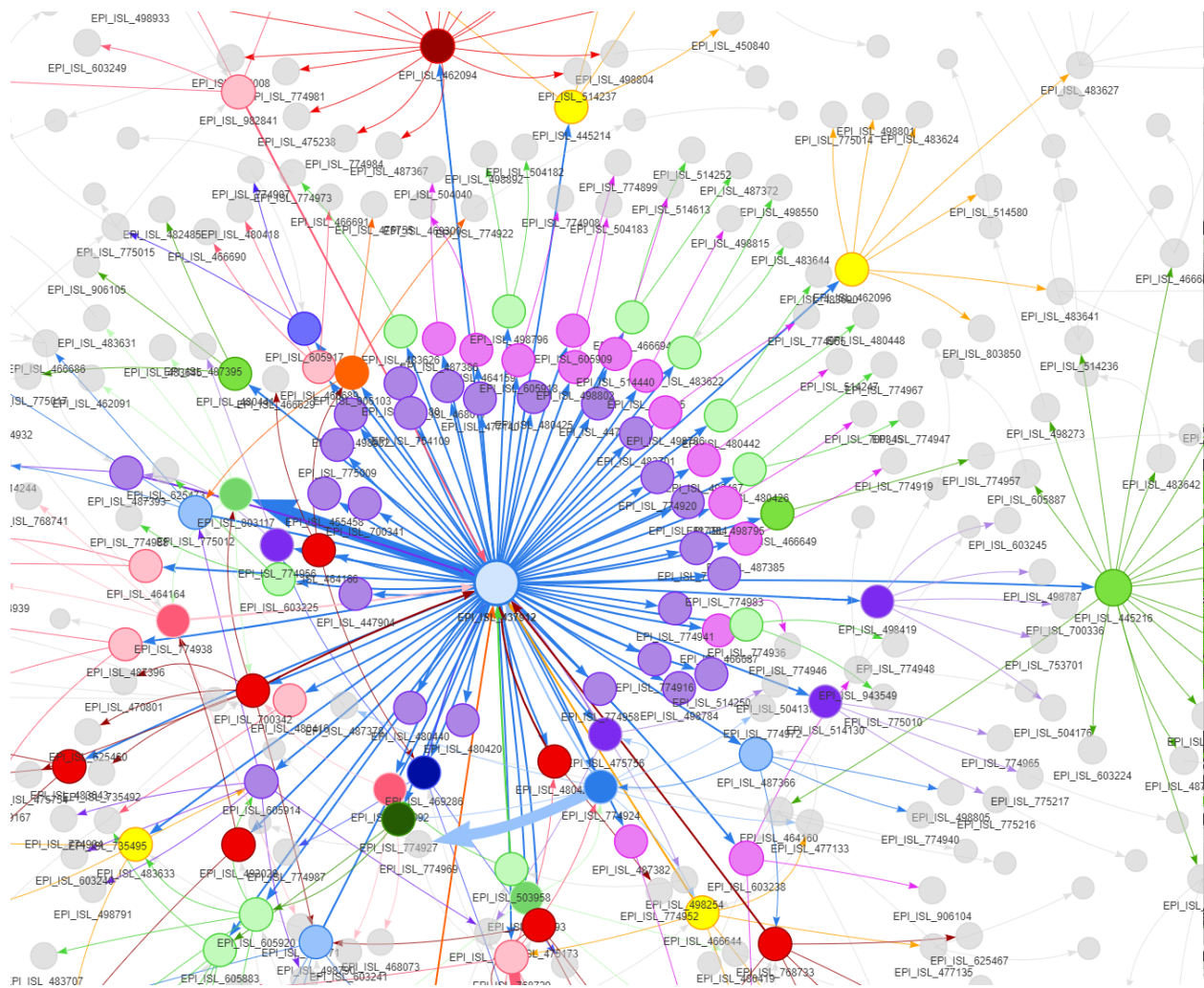
177   Dhaka was the hub of viral spread (Fig 2). Area-specific detailed chronological distribution of

178   SARS-CoV-2 variants is provided in the supplementary file (Supplementary file 2).

179



180

181   **Fig 2: Transmission of SARS-CoV-2 in Bangladesh.** Unlike others, Dhaka was connected with

182   all other parts of the country, therefore recognized as a viral transmission hub. Arrows at the tip of

183   the line dictate the direction of transmission.

184    To get a clearer idea of the viral circulation trend and back and forth transmission in different

185    divisions of the country, we extensively analyzed the variants chronologically. We figured out that

186    the whole country was mostly filled with a few major lineages throughout the times, but

187    interestingly their dominance varied. We have seen that some lineages were missing from a

188    particular area at a particular time and then returned, maybe due to mass people's movement from

189    other areas. For example, B.1.1 lineages were present in Mymensingh from the very beginning till

190    June 2020. Then, this variant was missing there for five months but reappeared in the middle of

191    December 2020. However, the variant was found during this period in Dhaka and Chattogram. On

192    the other hand, the sub-lineages of Beta variant B.1.351.3 were missing in Sylhet for two months

193    from February to March 2021 and occurred again in April 2021, while this variant was present in

194    other divisions during this time. Several other back and forth circulation of strains were observed,

195    for example, AY.100 and AY.102 in Dhaka. Detailed circulation of the variants information is

196    provided in the supplementary file (Supplementary file 2).

197

198



**Fig 3: Strain to strain transmission network of the SARS-CoV-2 in Bangladesh.** The sizes of nodes are proportional to the number of sequences a cluster contains and the thickness of the lines and arrows represent the frequency of transmission. The arrows reflect the direction of transmission among the viral clusters. The first sequence from the country is at the center of the network, and different clusters originated from the very first sequence, which gave rise to further subgroups; eventually, tips of the network reached.

Finally, we have built a viral transmission network using all our analysis data set sequences. Dhaka was found to be the center of viral transmission and directly connected with all other locations,

208 while others were not. For example, we did not find any direct connection between Chattogram

209 with Khulna and Mymensingh, Rangpur with Khulna, and Barisal did not have any connection

210 with Sylhet (Fig 2). In addition, a strain-specific transmission network reveals the connections

211 among different clusters and routes of viral spread from root to tip (Fig 3). With the time-calibrated

212 analysis, we have observed that the sequences from Dhaka remain at the center of the network and

213 determine the course of transmission forming connections with several subgroups.

214

215 **Mutation analysis summary**

216 Till the present study, we have found 7659 unique mutations present in 4622 sequences where 482

217 were extragenic mutations, and the rest were in the coding regions. In the coding region, a total of

218 4103 missense, 2865 synonymous, ten insertion, 125 deletion and 74 premature stop codon

219 mutations were observed (Fig 4A). Moreover, our analysis demonstrated 37.64 mutations per

220 sequence, where 24.61 mutations were missense, and the ratio of acquiring missense over

221 synonymous mutations increased gradually (Fig 4B). We have seen the number of mutations

222 increased gradually over time, yet nearly 29% of the sequences carried mutations below 30, and

223 more than 55.25% of sequences had 30 to 50 mutations. The highest number of mutations detected

224 was 78 in two strains isolated from Dhaka on 28[th] February 2022, and the lowest number was only

225 one found in a sequence from 11[th] May 2021. Fig 4B clearly demonstrates two remarkable rises in

226 mutations, one in February 2021 due to the introduction of Delta variants. Another sharp rise was

227 observed in January 2022 because of the highly transmissible Omicron variant with a large number

228 of mutations in the spike protein. However, the individual genes went through mutation

229 distinctively. Therefore, we thoroughly carried out the mutational analysis of all the SARS-CoV-

230 2 sequences from Bangladesh and summarized the results in table 1 and figure 4.

231 **Table 1: SARS-CoV-2 mutation summary on individual genes.**

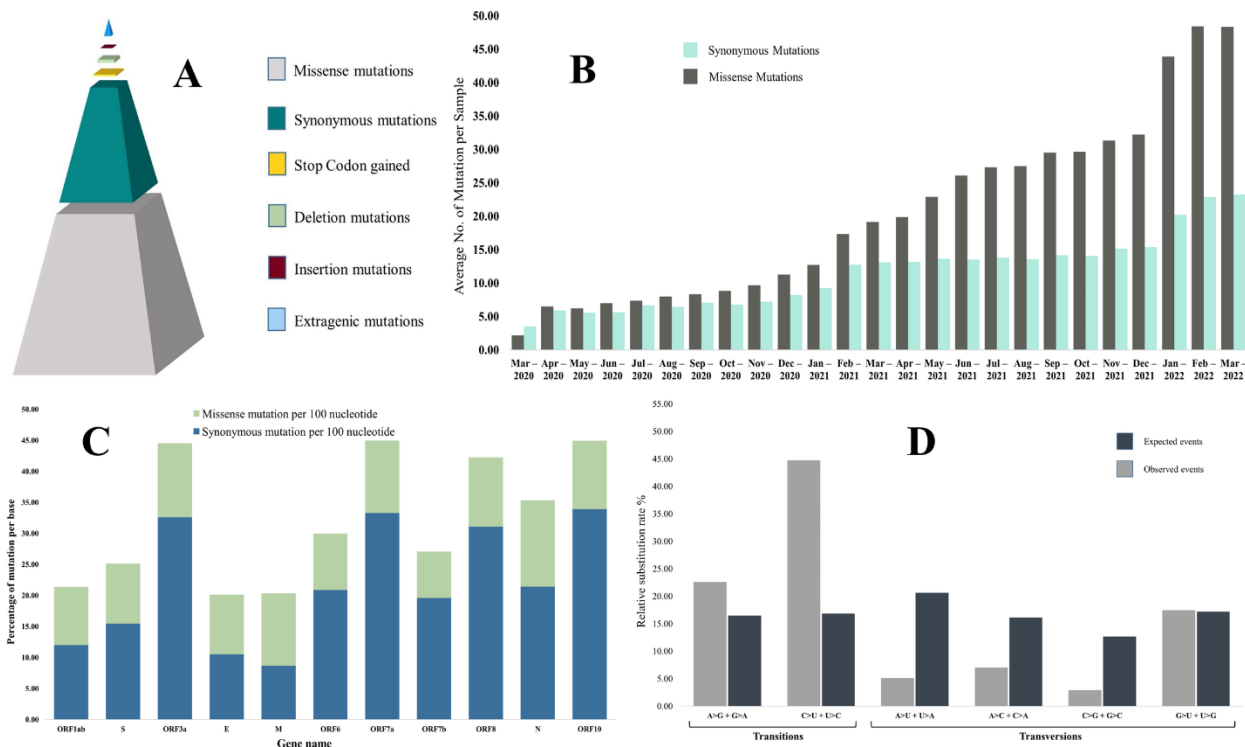| ORF | No. of non-mutant sequences | Percentage of mutated sequences | No. of synonymous mutations | No. of missense mutations | Percentage of missense mutations | Mutation density | No. of frequent mutations (n >= 10%) | No. of insertion mutation | No. of deletion mutation | No. of stop codon gained |
|---|---|---|---|---|---|---|---|---|---|---|
| ORF1ab | 1 | 99.98% | 1997 | 2550 | 56.08% | 21.36% | 29 | 2 | 33 | 20 |
| S | 2 | 99.96% | 370 | 590 | 61.46% | 25.12% | 31 | 4 | 40 | 11 |
| ORF3a | 873 | 81.11% | 99 | 270 | 73.17% | 44.57% | 4 | 2 | 5 | 1 |
| E | 3395 | 26.55% | 22 | 24 | 52.17% | 20.18% | 1 | 0 | 0 | 1 |
| M | 1423 | 69.21% | 78 | 58 | 42.65% | 20.33% | 3 | 0 | 1 | 2 |
| ORF6 | 3838 | 16.96% | 17 | 39 | 69.64% | 30.11% | 1 | 0 | 4 | 3 |
| ORF7a | 2258 | 51.15% | 43 | 122 | 73.94% | 45.08% | 3 | 0 | 11 | 11 |
| ORF7b | 2395 | 48.18% | 10 | 26 | 72.22% | 27.27% | 2 | 0 | 4 | 4 |
| ORF8 | 1604 | 65.30% | 41 | 114 | 73.55% | 42.35% | 1 | 1 | 11 | 13 |
| N | 78 | 98.31% | 175 | 270 | 60.67% | 35.32% | 12 | 1 | 15 | 3 |
| ORF10 | 4250 | 8.05% | 13 | 40 | 75.47% | 45.30% | 0 | 0 | 1 | 5 |

232

233    ORF10 and ORF7a harbored the highest mutation density with 45.30% and 45.08% mutations per

234    base, respectively, although only 8.05% of sequences were found to carry mutations in ORF10.

13

235    On the other hand, 99.98% and 99.96% of sequences had mutations in ORF1ab and S genes, but

236    their mutation density was lower at 21.36% and 25.21%, respectively. ORF6 was found to be the

237    most stable gene of SARS-CoV-2 in sequences from Bangladesh, with only 16.96% sequences

238    carrying the mutations, 30.11%% mutations per base and 69.64% missense mutations. ORF3a was

239    identified to harbor the highest percentage (75.69%) of missense mutations. In comparison, the

240    least percentage of missense mutations (42.65%) with 22.33% mutations per base was found in

241    membrane protein-encoding gene M. It was clearly evident that non-structural proteins were

242    subjected to more missense mutations than non-synonymous mutations compared with structural

243    proteins (Fig 4C). In addition, we have found several deletions and insertion mutations where both

244    the highest occurrences were found in the spike protein-coding S gene with 40 unique deletions

245    and four insertions (Table 1). On the other hand, the highest number of unique stop codons were

246    present in ORF1ab, with 40 out of 74 total stop codon mutations detected (table 1).

247    Among the 7786 mutations, 6968 were SNP, where 4697 and 2271 were involved in transition and

248    transversion events, respectively, rendering a transition transversion ratio of 2.07. Transition

249    mutations were calculated to be more prevalent than expected if mutational events took place

250    randomly, which clearly revealed the nucleotide substitution bias (Fig 4D). Then, transition

251    mutation C>U was the most frequent event, being 30.67% of total mutations and 45.50% of

252    transition mutations, followed by the transversion event G>U, which was 15.37% of the total

253    mutations (Supplementary file 3).

254

**Fig 4: Summary of mutational events. A.** Type of mutations among the sequences. Considering all the unique mutations, missense mutations were found to be the most prevalent event. **B.** The average number of mutations per sequence in each month. Although the number of mutations gradually increased with time, we observed a sharp increase from January 2021 when the Alpha variant entered the country. Moreover, more non-synonymous mutations emerged with time than synonymous mutations. **C.** Percentage of mutation per base in each gene. ORF3a had the highest density of mutations, while Envelop protein is the least mutated. Missense mutations are more prevalent than synonymous mutations. **D.** Nucleotide substitution rates for each of the four nucleotides among the SARS-CoV-2 genomes. Transition events were more prevalent than transversion events. C>U substitution rate was more than three times higher than the expected rate.

Then, out of the ten most prevalent mutations in Bangladesh, three were extragenic, one was synonymous, and six were missense mutations, where 23403A>G (missense mutation) was the

15

269  highest prevalent, followed by the second highest 14408C>T (missense mutation) which resembles

270  the global scenario and these two mutations were accompanied by 3037>C>T (synonymous

271  mutation). Among the top 7 mutations in the coding region, three were in the spike protein

272  (D614G, P681R and T478K), two were in the ORF1ab (P4715L and F924F), one was in membrane

273  protein (I82T), and another was in ORF3a (S26L). These seven mutations had a high prevalence

274  globally because these belong to different variants of concerns. In addition, 1163A>T (nsp2:

275  I120F) was a highly prevalent and unique mutation found in Bangladesh from the beginning of the

276  pandemic while it was absent in other countries. This mutation was present in more than 21% of

277  sequences. Interestingly enough, from linkage disequilibrium (LD) analysis, we have found that

278  all the mentioned mutations had a very strong correlation ($R^2 = 1.00$) and occurred in parallel since

279  their first appearance (Fig 5). 1163A>T mutation was predicted to have been 100% ($R^2 = 1.00$)

280  connected with 20 more mutations in parallel, considering the high frequent mutations present at

281  least in 10% of our sequence set. Additionally, several other mutations were occurring in parallel

282  with very strong LD values due to the introduction of several variants of concerns in the country

283  (Fig 5).

284

**Fig 5: Linkage disequilibrium plot.** The LD plot is generated considering the most prevalent SNPs. The number at the top denotes the SNP position, and squares are colored by standard (D'/LOD). The brighter red color indicates a higher D′ value and vice versa. The number in square is $r^2$ value.

**Effect of the mutations**

The mutations affect viral infectivity, transmissibility, virulence, viral fitness, selection pressure, proteome structure, and evolution. Overall, the SARS-CoV-2 genomes had very high nucleotide identity with an average of 37.64 mutations and low overall nucleotide diversity ($\pi$) 0.004. Although overall nucleotide diversity was lower, it varied from gene to gene. For example, ORF8 had the highest nucleotide diversity (0.01543), while gene ORF10 was most stable with a $\pi$ value of 0.00059 (Table 2). Analyzing the Bangladeshi sequences, the most diverse spot of the genome was in the spike protein gene at position 23009 with nucleotide diversity value $\pi$=0.16372 while the least diverse spot found was at position 11069 of ORF1ab with a $\pi$ value of 0.00005. Nevertheless, most of the genes overall had lower nucleotide diversity, which signifies selective

17

301    sweep due to increased mutations that benefit the strains and lead to the reduction of genetic

302    variation. In addition, we have found that three genes (ORF3a, ORF7b and ORF10) were under

303    positive selection pressure or directional selection because the mutations present in them were

304    advantageous to them; therefore, their frequencies were on the rise while the rest of the genes were

305    facing negative evolution pressure to stabilize against the deleterious mutations they have got from

306    random mutational events. Precisely, only 47 sites were facing positive or diversifying selection

307    pressure against 190 sites found to be under negative or purifying pressure to stabilize the genomic

308    variations (Supplementary file 3).

309

310    **Table 2: Summary of Mutation's effect on each protein.**

| ORF | Nucleotide diversity ($\pi$) | dN/dS | No. of sites towards positive selection | No. of sites towards negative selection |
|---|---|---|---|---|
| ORF1ab | 0.0017 | 0.579 | 28 | 103 |
| S | 0.00526 | 0.74 | 7 | 48 |
| ORF3a | 0.00292 | 1.479 | 2 | 11 |
| E | 0.00223 | 0.479 | 0 | 2 |
| M | 0.00206 | 0.371 | 1 | 6 |
| ORF6 | 0.00452 | 0.928 | 5 | 18 |
| ORF7a | 0.00549 | .987 | 0 | 3 |

18

| ORF7b | 0.0046 | 1.44 | 0 | 0 |
|---|---|---|---|---|
| ORF8 | 0.01543 | 0.941 | 1 | 9 |
| N | 0.00579 | 0.841 | 3 | 14 |
| ORF10 | 0.00059 | 1.17 | 0 | 2 |

311  *(dN/dS> 1 is positive selection, dN/dS = 1 neutral selection, dN/dS< 1 negative selection, dN/dS = 0 is conserved*

312  *region)*

313

314



316  **Fig 6 Distribution of missense mutations in the genome.** The height of the spikes is proportional

317  to the number of sequences that got mutation at that location. Regions between these spikes are

318  stable, which could be targeted for further vaccine and therapeutics development.

319

320  Finally, these mutations affected the virus from the evolutionary perspective and shook the

321  stability of the proteins they encode. Most of the mutations were previously reported to affect the

322  stability of the whole proteome of SARS-CoV-2 negatively. However, all the genes were not

19

323    affected to the same extent by mutational events (Fig 6). For example, only 42.65% of mutations

324    on the membrane protein-coding M gene were missense which was 73.55% in the case of ORF3a

325    (table 1). As of now, vaccines and therapies target the spike protein, which is highly mutated.

326    That's one reason why people continue to develop symptoms after successful vaccination. It is

327    possible that current vaccines and therapies will not work in future due to a high number of

328    mutations occurring. The less affected genes could therefore be targeted for medicine and vaccine

329    development. Figure 6 shows spikes that represent mutations, and the height of the spikes is

330    proportional to the number of mutations that have taken place at that position in the genome. As

331    we can see, there are plenty of stable regions between the spikes, which could be targeted for

332    therapeutics and vaccine development against SARS-CoV-2.

333

334    **Discussion**

335    SARS-CoV-2 has been circulating in Bangladesh for over two years, and many strains are

336    sequenced from different parts of the country, helping us carry out downstream analysis to depict

337    different variants, transmission, and evolution inside the country. Investigating 4622 whole-

338    genome sequences from Bangladesh, we have seen B.1.1.25 lineage was in dominance since the

339    beginning along with other lineages containing a small number of sequences, but since March

340    2021, strains from B.1.1.25 lineage are surmounted by another lineage B.1.351.3, which is a sub-

341    line of Beta variant. In addition to that, we have observed a slight increase of Alpha variants for a

342    short time since their first appearance in January 2021. However, In April 2021 Delta variant

343    emerged and dominated other variants until the arrival of Omicron. The Omicron variant

344    comprises three main sub-lineages termed BA.1, BA.2 and BA.3. Both BA.1 and BA.2 are found

345    in Bangladesh, and currently, BA.2 is the dominant variant. Although BA.1 and BA.2 have

20

346     numerous mutations in common, 20 mutations in the spike protein differentiate the two sub-

347     lineages, and BA.2 also displays a marked decreased sensitivity to many neutralizing monoclonal

348     antibodies (mAbs) when compared to previous VOCs (37). Therefore, with further mutations, this

349     BA.2 sub-lineage is keeping the risk of having another COVID-19 wave alive in the country.

350     On the other hand, geographical analysis depicts Dhaka and Chattogram containing a more

351     diversified number of sequences than other parts of the country. Our analysis has limitations at

352     this point because we had a higher number of sequences from these two regions than others. The

353     sequences were more diversified in the first phase of the pandemic. However, with the arrival of

354     the Delta and Omicron variants, the divergence reduced drastically, maybe due to viral adaptation

355     following the "Survival of the fittest" theory of natural selection, although we have found several

356     sub-lineages of the Delta variant. Additionally, we have seen Dhaka being the viral transmission

357     hub, which is obvious since it is the capital city of Bangladesh, but this city is not the only

358     transmission source. From extensive analysis, we have built the SARS-CoV-2 transmission

359     network between different administrative divisions and observed the back and forth transmission

360     of the virus inside Bangladesh. This situation arose due to a lack of restriction on the mass

361     movement; public gatherings were not limited duly, and other socioeconomic events.

362     From the mutational perspective, we have seen a total of 7659 unique mutations present in 4622

363     sequences with 37.64 mutations per sample where on average 24.61were coding variants, which

364     happens to be significantly higher than the global average of 7.23, reported in July 2020 (38). This

365     sharp rise of mutations indicates the SARS-CoV-2 might be facing strong challenges from the

366     host's immunologic response in addition to random regular mutational events of RNA viruses,

367     which is one of the reasons for the emergence of new variants of concerns. At the nucleotide level,

368     67.41% of the mutations were transition events, and a molecular bias was present for C>U, which

369    tends to mutate hydrophilic amino acids into hydrophobic ones (36). Overall, 1109 unique C>U

370    missense mutation on the genome was observed, where 194, 289, and 162 of the mutations were

371    skewed towards phenylalanine, isoleucine, and leucine codons, respectively. Moreover, 236 C>U

372    mutations were involved in altering the proline, which is known to be a strong helix breaker (39).

373    Therefore, proline to another amino acid shift might have a deleterious effect on the  SARS-CoV-

374    2 proteome.

375    On the other hand, we have found that most of the genes were under negative selection pressure

376    while only three non-structural protein-coding genes were under Darwinian (positive) selection,

377    which indicates that most of the random mutational events were deleterious for the SARS-CoV-2

378    (40), maybe due to the immunologic potential of people of Bangladesh and our demography.

379    However, the dN/dS ratio of the receptor-binding region (RBD) of the spike protein was higher,

380    indicating that mutations in this region were advantageous. The result correlates with the

381    emergence of different variants of concerns like Alpha, Beta, Delta and Omicron. The RBD region

382    is considered the most important part of the virus since it attaches to ACE2 during viral infection

383    to host cells. Those advantageous mutations might increase pathogenicity, infectivity,

384    transmissibility, and letting it evade the host immunity (41–43). Moreover, most of the current

385    therapies and vaccines are developed targeting the BRD-ACE2 interaction. Therefore, a higher

386    dN/dS ratio also warns us about the emergence of new deadly variants in future with further

387    mutations in this region and vaccine failure. However, analyzing all the sequences from

388    Bangladesh, we have seen that the whole proteome was not affected to the same extent. There were

389    regions of high and low nucleotide diversity. While highly affected regions are evolving faster,

390    regions with low nucleotide divergence would render us the opportunity to develop new vaccines

391 and antibodies for treatment and development detection kits to reduce the number of false test

392 reports during this pandemic.

393 To sum up, considering the limitations regarding sequence number variations in different parts of

394 the country, we have analyzed all the unique and global mutations present in Bangladesh, which

395 are thoroughly reported in the supplementary files. This data would facilitate researchers further

396 from various perspectives like investigating viral transmission, the connection among isolates,

397 evolution patterns, and dynamics of divergence of the virus.

398

## Methods and materials

### Sequence retrieval and Lineage determination

401 Using completeness and coverage filters on the sequences, all the SARS-CoV-2 genomes

402 submitted from Bangladesh until 25th March 2022 were retrieved from the Global Initiative on

403 Sharing All Influenza Data (GISAID) database (www.gisaid.org) (9). Prior to downstream

404 analysis, all sequences were quality checked and sequences with more than 5% ambiguous

405 characters were omitted. The sorted sequences were then classified by Phylogenetic Assignment

406 of Named Global Outbreak LINeages (Pangolin) with COVID-19 Lineage Assigner

407 (https://pangolin.cog-uk.io/) (20). Furthermore, we excluded lineages carrying less than ten

408 sequences to address the important lineages. We analyzed and visualized the sequence lineage

409 distribution in R within the country.

### Transmission analysis:

411 First, the selected sequences were aligned using the Mafft algorithm (21), followed by the

412 construction of a maximum likelihood phylogenetic tree using IQ-TREE (22)and calibrating the

413   tree based on time with TreeTime (23). Using the StrainHub tool (24), we built the SARS-CoV-2

414   transmission network in Bangladesh from the reconstructed tree and metadata.

415   **Mutation Analysis:**

416   We have aligned each sequence with the reference sequence (NC_045512.2) (8) using the

417   minimap2 algorithm (25) and called the variants with Samtools (26). Additionally, SNP-sites (27),

418   CovSeq (28) and an online server Coronapp (29) were used to detect the mutations present in the

419   sequences and the common mutations from these four sources were considered. Finally, SNPeff

420   was used to predict the impact of the mutations (30).

421   **Effects of mutation:**

422   First of all, we used TASSEL software (31) to determine the nucleotide diversity ($\pi$) using a 20

423   base-pair window at five base-pair steps. Then we calculated the direction of selection in the

424   sequences to know if diversity moves away from neutrality and to understand the pattern of

425   evolution using the SLAC algorithm (32) in the HyPhy software package (33). Linkage

426   disequilibrium among mutations prevalent in 10% or more sequences were calculated using

427   AutoVem (34) and presented by the R2 index using HaploView (35). Then, along with determining

428   the nucleotide substitution bias, the expected and observed transition, transversion events as well

429   as their ratio were calculated by the method used by Matyášek R, Kovařík A (36).

430

431   **Acknowledgements**

459     Yuan ML, Zhang YL, Dai FH, Liu Y, Wang QM, Zheng JJ, Xu L, Holmes EC, Zhang YZ.

460     2020. A new coronavirus associated with human respiratory disease in China. Nature

461     579:265–269.

462   9.   Elbe S, Buckland-Merrett G. 2017. Data, disease and diplomacy: GISAID's innovative

463     contribution to global health. Glob Challenges 1:33–46.

464   10.   Worldometer. 2021. COVID Live Update: 239,169,612 Cases and 4,875,781 Deaths from

465     the Coronavirus. Worldometer.

466   11.   Islam MT, Talukder AK, Siddiqui MN, Islam T. 2020. Tackling the COVID-19 pandemic:

467     The Bangladesh perspective. J Public health Res 9:389–397.

468   12.   Saha S, Malaker R, Sajib MSI, Hasanuzzaman M, Rahman H, Ahmed ZB, Islam MS,

469     Islam M, Hooda Y, Ahyong V, Vanaerschot M, Batson J, Hao S, Kamm J, Kistler A, Tato

470     CM, DeRisi JL, Saha SK. 2020. Complete Genome Sequence of a Novel Coronavirus

471     (SARS-CoV-2) Isolate from Bangladesh. Microbiol Resour Announc 9.

472   13.   Choi JY, Smith DM. 2021. SARS-CoV-2 variants of concern. Yonsei Med J

473     https://doi.org/10.3349/ymj.2021.62.11.961.

474   14.   Sanyaolu A, Okorie C, Marinkovic A, Haider N, Abbasi AF, Jaferi U, Prakash S,

475     Balendra V. 2021. The emerging SARS-CoV-2 variants of concern. Ther Adv Infect Dis

476     https://doi.org/10.1177/20499361211024372.

477   15.   Ogando NS, Zevenhoven-Dobbe JC, van der Meer Y, Bredenbeek PJ, Posthuma CC,

478     Snijder EJ. 2020. The Enzymatic Activity of the nsp14 Exoribonuclease Is Critical for

479     Replication of MERS-CoV and SARS-CoV-2. J Virol 94.

480   16.   Gribble J, Stevens LJ, Agostini ML, Anderson-Daniels J, Chappell JD, Lu X, Pruijssers

481     AJ, Routh AL, Denison MR. 2021. The coronavirus proofreading exoribonuclease

482        mediates extensive viral recombination. PLoS Pathog 17.

483   17.   Duchêne S, Ho SY, Holmes EC. 2015. Declining transition/transversion ratios through

484        time reveal limitations to the accuracy of nucleotide substitution models. BMC Evol Biol

485        15.

486   18.   Shishir TA, Naser I Bin, Faruque SM. 2021. In silico comparative genomics of SARS-

487        CoV-2 to determine the source and diversity of the pathogen in Bangladesh. PLoS One 16.

488   19.   Rahman MM, Kader SB, Rizvi SMS. 2021. Molecular characterization of SARS-CoV-2

489        from Bangladesh: implications in genetic diversity, possible origin of the virus, and

490        functional significance of the mutations. Heliyon 7:e07866.

491   20.   O'Toole Á, Scher E, Underwood A, Jackson B, Hill V, McCrone JT, Colquhoun R, Ruis

492        C, Abu-Dahab K, Taylor B, Yeats C, du Plessis L, Maloney D, Medd N, Attwood SW,

493        Aanensen DM, Holmes EC, Pybus OG, Rambaut A. 2021. Assignment of epidemiological

494        lineages in an emerging pandemic using the pangolin tool. Virus Evol 7.

495   21.   Katoh K, Rozewicki J, Yamada KD. 2018. MAFFT online service: Multiple sequence

496        alignment, interactive sequence choice and visualization. Brief Bioinform 20:1160–1166.

497   22.   Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, Von Haeseler A,

498        Lanfear R, Teeling E. 2020. IQ-TREE 2: New Models and Efficient Methods for

499        Phylogenetic Inference in the Genomic Era. Mol Biol Evol 37:1530–1534.

500   23.   Sagulenko P, Puller V, Neher RA. 2018. TreeTime: Maximum-likelihood phylodynamic

501        analysis. Virus Evol 4.

502   24.   De Bernardi Schneider A, Ford CT, Hostager R, Williams J, Cioce M, Çatalyürek Ü V.,

503        Wertheim JO, Janies D. 2020. StrainHub: A phylogenetic tool to construct pathogen

504        transmission networks. Bioinformatics 36:945–947.

505    25.    Li H. 2018. Minimap2: Pairwise alignment for nucleotide sequences. Bioinformatics

506             34:3094–3100.

507    26.    Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G,

508             Durbin R. 2009. The Sequence Alignment/Map format and SAMtools. Bioinformatics

509             25:2078–2079.

510    27.    Page AJ, Taylor B, Delaney AJ, Soares J, Seemann T, Keane JA, Harris SR. 2016. SNP-

511             sites: rapid efficient extraction of SNPs from multi-FASTA alignments. Microb genomics

512             2:e000056.

513    28.    Simonetti M, Zhang N, Harbers L, Milia MG, Brossa S, Huong Nguyen TT, Cerutti F,

514             Berrino E, Sapino A, Bienko M, Sottile A, Ghisetti V, Crosetto N. 2021. COVseq is a

515             cost-effective workflow for mass-scale SARS-CoV-2 genomic surveillance. Nat Commun

516             12.

517    29.    Mercatelli D, Triboli L, Fornasari E, Ray F, Giorgi FM. 2021. Coronapp: A web

518             application to annotate and monitor SARS-CoV-2 mutations. J Med Virol 93:3238–3245.

519    30.    Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden

520             DM. 2012. A program for annotating and predicting the effects of single nucleotide

521             polymorphisms, SnpEff: SNPs in the genome of Drosophila melanogaster strain w1118;

522             iso-2; iso-3. Fly (Austin) 6:80–92.

523    31.    Bradbury PJ, Zhang Z, Kroon DE, Casstevens TM, Ramdoss Y, Buckler ES. 2007.

524             TASSEL: Software for association mapping of complex traits in diverse samples.

525             Bioinformatics 23:2633–2635.

526    32.    Kosakovsky Pond SL, Frost SDW. 2005. Not so different after all: A comparison of

527             methods for detecting amino acid sites under selection. Mol Biol Evol 22:1208–1222.

528    33.    Kosakovsky Pond SL, Poon AFY, Velazquez R, Weaver S, Hepler NL, Murrell B, Shank

529            SD, Magalis BR, Bouvier D, Nekrutenko A, Wisotsky S, Spielman SJ, Frost SDW, Muse

530            S V. 2020. HyPhy 2.5 - A Customizable Platform for Evolutionary Hypothesis Testing

531            Using Phylogenies. Mol Biol Evol 37:295–299.

532    34.    Xi B, Jiang D, Li S, Lon JR, Bai Y, Lin S, Hu M, Meng Y, Qu Y, Huang Y, Liu W,

533            Huang L, Du H. 2021. AutoVEM: An automated tool to real-time monitor epidemic trends

534            and key mutations in SARS-CoV-2 evolution. Comput Struct Biotechnol J 19:1976–1985.

535    35.    Barrett JC, Fry B, Maller J, Daly MJ. 2005. Haploview: Analysis and visualization of LD

536            and haplotype maps. Bioinformatics 21:263–265.

537    36.    Matyášek R, Kovařík A. 2020. Mutation patterns of human SARS-CoV-2 and bat

538            RATG13 coronavirus genomes are strongly biased towards C>U transitions, indicating

539            rapid evolution in their hosts. Genes (Basel) 11:1–13.

540    37.    Planas, D., Saunders, N., Maes, P. et al. Considerable escape of SARS-CoV-2 Omicron to

541            antibody neutralization. Nature 602:671–675

542    38.    Mercatelli D, Giorgi FM. 2020. Geographic and Genomic Distribution of SARS-CoV-2

543            Mutations. Front Microbiol 11.

544    39.    Li SC, Goto NK, Williams KA, Deber CM. 1996. α-Helical, but not β-sheet, propensity of

545            proline is determined by peptide environment. Proc Natl Acad Sci U S A 93:6676–6681.

546    40.    Lin JJ, Bhattacharjee MJ, Yu CP, Tseng YY, Li WH. 2019. Many human RNA viruses

547            show extraordinarily stringent selective constraints on protein evolution. Proc Natl Acad

548            Sci U S A 116:19009–19018.

549    41.    Lan J, Ge J, Yu J, Shan S, Zhou H, Fan S, Zhang Q, Shi X, Wang Q, Zhang L, Wang X.

550            2020. Structure of the SARS-CoV-2 spike receptor-binding domain bound to the ACE2

551        receptor. Nature 581:215–220.

552   42.   Barros EP, Casalino L, Gaieb Z, Dommer AC, Wang Y, Fallon L, Raguette L, Belfon K,

553        Simmerling C, Amaro RE. 2021. The flexibility of ACE2 in the context of SARS-CoV-2

554        infection. Biophys J 120:1072–1084.

555   43.   Xu C, Wang Y, Liu C, Zhang C, Han W, Hong X, Wang Y, Hong Q, Wang S, Zhao Q,

556        Wang Y, Yang Y, Chen K, Zheng W, Kong L, Wang F, Zuo Q, Huang Z, Cong Y. 2021.

557        Conformational dynamics of SARS-CoV-2 trimeric spike glycoprotein in complex with

558        receptor ACE2 revealed by cryo-EM. Sci Adv 7.

559