

1 **Tajima *D* test accurately forecasts Omicron / COVID-19 outbreak**

2

3 **Ting-Yu Yeh and Gregory P. Contreras**

4

5 Agricultural Biotechnology Laboratory, Auxergen Inc., Columbus Center, 701 E Pratt

6 Street, Baltimore, MD 21202, USA

7

8 *Ting-Yu Yeh is Corresponding author

9

10 E-mail: yehty@auxergen.com

11 Tel: 1-443-762-1974

12

13

14 Key words: COVID-19, SARS-CoV-2, Omicron variants, polymorphism, site-frequency

15 spectrum, Tajima *D*

16 Short running title: Forecasting SARS-CoV-2 Omicron outbreak

17 Word count: 1893 words (text)

18

19

20 **Abstract**

21

22 On 26 November 2021, the World Health Organization designated the SARS-CoV-2
23 variant B.1.1.529, Omicron, a variant of concern. However, the phylogenetic and
24 evolutionary dynamics of this variant remain unclear. An analysis of the 131
25 Omicron variant sequences from November 9 to November 28, 2021 reveals that
26 variants have diverged into at least 6 major subgroups. 86.3% of the cases have an
27 insertion at amino acid 214 (INS214EPE) of the spike protein. Neutrality analysis of
28 DH ($-2.814, p < 0.001$) and Zeng's E ($0.0583, p = 1.0$) tests suggested that directional
29 selection was the major driving force of Omicron variant evolution. The
30 synonymous (D_{syn}) and nonsynonymous (D_{nonsyn}) polymorphisms of the Omicron
31 variant spike gene were estimated with Tajima's D statistic to eliminate
32 homogenous effects. Both D ratio ($D_{nonsyn}/D_{syn}, 1.57$) and ΔD ($D_{syn} - D_{nonsyn}, 0.63$)
33 indicate that purifying selection operates at present. The low nucleotide diversity
34 (0.00008) and Tajima D value ($-2.709, p < 0.001$) also confirms that Omicron variants
35 had already spread in human population for more than the 6 weeks than has been
36 reported. These results, along with our previous analysis of Delta and Lambda
37 variants, also supports the validity of the Tajima's D test score, with a threshold
38 value as -2.50 , as an accurate predictor of new COVID-19 outbreaks.

39

40

41

42 **Introduction**

43

44 On 26 November 2021, World Health Organization designated the SARS-CoV-2
45 variant B.1.1.529, Omicron, a variant of concern based on its unique mutations, and
46 unusual features suggesting that therapeutic monoclonal antibodies, such as
47 Regeneron, may be less effective against Omicron. Evidences that new mutations in
48 Omicron could have an impact on viral transmission or the severity of illness are
49 still under investigation (World Health Organization, 2021). Omicron variants were
50 first reported in the Gauteng province, South Africa on November 9, 2021, and
51 shortly it was detected in recent travelers to Belgium, Botswana, Canada, Hong
52 Kong, Austria, Australia, Portugal, Israel, United Kingdom, Netherland, and the
53 United States. One of the major concerns is that Omicron variants contain more than
54 30 mutations to the spike protein (A67V, Δ69-70, T95I, G142D/Δ143-145,
55 Δ211/L212I, ins214EPE, G339D, S371L, S373P, S375F, K417N, N440K, G446S,
56 S477N, T478K, E484A, Q493R, G496S, Q498R, N501Y, Y505H, T547K, D614G,
57 H655Y, N679K, P681H, N764K, D796Y, N856K, Q954H, N969K, L981F). Some of
58 these changes have been previously identified in Delta or Alpha variants and are
59 linked to heightened infectivity and the ability to evade infection-blocking
60 antibodies (Callaway, 2021). The likelihood of higher transmission rates has led
61 multiple countries to respond quickly to the Omicron variant.

62

63 The rapid increase in Omicron variant cases was found in the Gauteng province in
64 November, particularly in schools and among young people. Preliminary evidence

65 from genotyping tests suggests that Omicron may have been in circulation for quite
66 some time in South Africa (Callaway, 2021). To date, phylogenetic and evolutionary
67 status of Omicron variants are still not clear. Here we analyze the phylogenetic
68 relationship and selection pressure among the 131 available sequences of Omicron
69 variants.

70

71 **Materials and Methods**

72

73 131 complete SARS-CoV-2 genome sequences excluding low coverage from
74 November 9 to November 28, 2021 in 9 countries: Austria ($N=1$), Australia ($N=1$),
75 Belgium ($N=1$), Canada ($N=1$), Botswana ($N=17$), Hong Kong, China ($N=2$), Italy
76 ($N=1$), South Africa ($N=105$), and United Kingdom (UK, $N=2$) were collected from
77 the Global Initiative on Sharing All Influenza Data (GISAID)
78 (<https://www.gisaid.org/>). Sequence data in this study is available and deposited at
79 Figshare website (10.6084/m9.figshare.17105090).

80

81 FASTA files of viral sequences were downloaded and first aligned using MAFFT 7
82 software (Kato and Standley, 2013, Kuraku et al, 2013). Phylogenetic relationships
83 between Omicron variants were analyzed using the neighbour-joining method and
84 Jukes-Cantor substitution model with bootstrap resampling number set as five. The
85 radial phylogenetic tree was generated by exporting the tree file in Newick format
86 by MAFFT. The FigTree software (version 1.4.2) was used to display the cladogram
87 (Rambaut, 2021).

88

89 The polymorphisms of the SARS-CoV-2 Omicron variants was analyzed based on the
90 site-frequency spectrum: (1) Tajima's D test, (2) normalized Wu and Fay's DH test,
91 and (3) Zeng's E test, using DNASP6 software with the bat SARS-CoV-2 WIV04
92 (GenBank MN996528.1) as the outgroup sequence (Rozas, et al., 2017). Statistical
93 significance of observed values of different tests was obtained by coalescent
94 simulation with intermediate recombination after 10,000 replicates. The probability
95 for each statistic was calculated as the frequency of replicates with a value lower
96 than the observed statistic (two-tailed test) by DNASP6 (Rozas, et al., 2017).

97

98 A modified Tajima's D statistic was used to examine purifying selection based on
99 non-synonymous ($D_{\text{non-syn}}$) versus synonymous sites (D_{syn}) of SARS-CoV-2 genes was
100 calculated as previously described (Hahn et al, 2002, Hughes, et al., 2005, Yeh and
101 Contreras, 2021). The average number of pairwise synonymous differences (kS)
102 and the average number of nonsynonymous pairwise differences (kN), the number
103 of synonymous segregating sites (S_S), and the number of nonsynonymous
104 segregating sites (S_N) were computed with equation: (1) $S^*_S = S_S/a_1$, (2) $S^*_N = S_N/a_1$,

105 where a_1 is defined as $a_1 = \sum_{i=1}^{n-1} \frac{1}{i}$ (Tajima, 1989, Hughes, et al., 2005). D_{syn} was
106 defined as $kS - S^*_S$, divided by the standard error of that difference, and $D_{\text{non-syn}}$ was
107 defined as $kN - S^*_N$, divided by the standard error of that difference (Tajima, 1989,
108 Hughes, et al., 2005).

109

110 **Results and Discussion**

111

112 **The insertion mutation Ins214EPE at the spike protein is present in some**

113 **Omicron variants**

114 The phylogenetic tree identifies at least 6 major subgroups of 131 Omicron variant

115 sequences after rooting with an outgroup virus sequence of SARS-CoV-2 WIV04

116 from Wuhan, China (Figure 1). 113 Omicron cases (86.3%) contain an insertion of

117 nine nucleotides (GAGCCAGAA) between nucleotide 22204 and 22205 according to

118 WIV04 sequence (Figure 2). This generates an insertion of three amino acids

119 (Glutamic acid-Proline-Glutamic acid) (INS214EPE) in the N-terminal domain (NTD)

120 of the spike protein.

121

122 Resende et al. has reported that most ins214 motifs were rare in sequences of

123 different lineages of SARS-CoV-2 (A.2.4, B.1, B.1.1.7, B.1.177, B.1.2, B.1.214, and

124 B.1.429) with an insertion motif of 3 or 4 amino acids (AKKN, KLGB, AQER, AAG,

125 KFH, KRI, and TDR) (Resende, et al, 2021). Interestingly, the ins214 with four amino

126 acids (ins214GATP, ins214GATP, ins214GATS) were also present in bat SC2r-CoV

127 isolated in China (RmYN02), Thailand (RacCS203), and Japan (Rc-o319),

128 respectively. Although amino acid identity at ins214 vary among SARS-CoV-2 or

129 SC2r-CoV lineages, the insertion size (3 or 4 amino acids) is also conserved in

130 Omicron. It suggested that this region of NTD is susceptible to a mutation or

131 insertion. Whether or not INS214EPE affects spike protein function or immune

132 response requires further investigation.

133

134 **Tajima's D test is able to forecast new Omicron outbreaks**

135 We have previously used neutrality tests to analyze the selection pressure of SARS-
136 CoV-2 in the shipboard quarantine on the Diamond Princess (Yeh and Contreras,
137 2021a) and the influence of full vaccination coverage on Delta variants among
138 different countries (Yeh and Contreras, 2021b). We also proposed that Tajima's D
139 test, a popular statistic in population and evolution genetics, can provide a
140 promising tool to forecast new COVID-19 outbreaks (Yeh and Contreras, 2021b).

141

142 Here the selection pressure of Omicron variants was first analyzed by multiple
143 neutrality tests. Within 131 Omicron variant sequences, 75 mutation sites were
144 identified with the nucleotide diversity (π , the average number of nucleotide
145 differences per site between two sequences) equal to 0.00008, which is significantly
146 lower than earlier outbreak of Delta variants in UK (0.0004-0.0006, $N=376$, March
147 26 to April 22, 2021), India (0.0004-0.0006, $N=67$, January 1 to March 11, 2021), or
148 Australia (0.0006, $N=75$, April 9 to May 6, 2021) (Yeh and Contreras, 2021b). Tajima
149 D test was calculated to compare π and total polymorphism (Tajima, 1989). Tajima's
150 D values were negative and significantly deviated from zero (-2.709 , $P < 0.001$)
151 among the whole genome sequences of Omicron variants. The negative values of
152 Tajima D were also detected in the ORF1ab (-2.617 , $P < 0.001$) and the spike gene ($-$
153 1.948 , $P < 0.005$). This result indicates an excess of nucleotide variants of low
154 frequency (Table 1), and strong selection and/or demographic expansion was
155 operating in Omicron.

156 We have previously shown that Tajima D values decreased as SARS-CoV-2 Delta and
157 lambda variants spread in human populations. One to three weeks after Tajima D
158 fell below -2.50, Delta variant outbreaks emerged in India and the UK (Yeh and
159 Contreras, 2021b). Taken together, the low π and Tajima D values suggested that
160 Omicron variants have most likely spread within a population weeks before the
161 samples were collected. **This data also agreed with our previous proposal that**
162 **Tajima's D test is useful to forecast new COVID-19 outbreaks regardless of the**
163 **sample sizes of different variants** (Yeh and Contreras, 2021b).

164

165 **Purifying selection was operating in Omicron variants**

166 Application of Tajima D is limited by the difficulty in distinguishing the influence of
167 both selection pressure and demographic expansion. To overcome this problem, we
168 included normalized DH test and Zeng's E test in our analysis (Zeng et al., 2006, Yeh
169 and Contreras, 2021b). The normalized DH test is affected by directional selection
170 but insensitive to demographic expansion. Zeng's E test is very sensitive to
171 population growth immediately after a sweep (Zeng et al., 2006). Genomic
172 polymorphisms of Omicron variants showed significantly negative DH values (-
173 2.814, $P < 0.05$), but E values were not significantly different from zero (0.0583,
174 $P > 0.1$) (Table 1). The results were similar after confining our analysis to the spike
175 gene (DH , -9.014, $p < 0.001$; Zeng's E , 5.504, $p < 0.001$) (Table 1), suggesting that
176 directional selection was the major driving force, without significant influence by
177 the demographic expansion.

178

179 We and others have shown the limitations to application of the inter-species
180 divergence, the dN/dS (ω) test, for analyzing selection pressure of SARS-CoV-2
181 previously (Mugal, et al., 2014, Kang et al., 2021, Yeh and Contreras, 2021a, Yeh and
182 Contreras 2021b). Therefore, here we determined the purifying selection using a
183 modified Tajima's D statistics instead of dN/dS (ω) test. Under purifying selection,
184 the frequency distribution of non-synonymous polymorphisms is negatively skewed
185 relative to the distribution of synonymous polymorphisms. Therefore, it takes more
186 negative values for non-synonymous ($D_{\text{non-syn}}$) than for synonymous sites (D_{syn}) of a
187 given gene (Hahn et al, 2002, Hughes, et al., 2005, Yeh and Contreras, 2021b). One
188 of the major advantages of $D_{\text{non-syn}}$ and D_{syn} analysis is that it is independent of
189 sample size, which allows us to compare $D_{\text{non-syn}}$ and D_{syn} values among different data
190 sets (Hughes et al., 2008). The same excess of low-frequency alleles in non-
191 synonymous polymorphism was also shown by the $D_{\text{non-syn}}$ values of the spike (-
192 1.771, $p < 0.001$) and N gene (-1.943, $p < 0.005$) (Table 1). **We conclude that**
193 **purifying selection led to constraints on the neutral mutations at non-**
194 **synonymous sites of the spike gene of Omicron variants.**

195 Negative values of Tajima D could be caused by a bottleneck event rather than
196 selection, but the bottleneck effect should affect all types of polymorphism equally
197 (Tajima, 1989). Hahn et al. have reported that $D_{\text{non-syn}}$ is disproportionately lower than
198 D_{syn} , because non-synonymous and synonymous mutations are affected unequally
199 by purifying selection (Hahn et al., 2002). The value of ΔD ($D_{\text{syn}} - D_{\text{non-syn}}$) increased as
200 purifying selection became stronger, and ΔD can eliminate the homogenous effects
201 (demographic expansion, selective sweep etc.). ΔD values of the spike and N gene of

202 SARS-CoV-2 Omicron variants are positive (0.643 and 0.377, $p < 0.005$) (Table 1).
203 Austin Hughes showed that the standard error terms in both equations of $D_{\text{non-syn}}$
204 and D_{syn} cancel out the sample size effect when D ratio ($D_{\text{non-syn}}/D_{\text{syn}}$) was applied to
205 the data analysis (Hughes, 2008). **D ratio values of the spike and N gene were**
206 **significantly more than 1 (1.57 and 1.24, Table 1), suggesting that purifying**
207 **selection pressure of Omicron spike and N gene was operating.** The values of
208 ΔD and D ratio of ORF1ab (-0.05 and 0.980, $p < 0.001$) were insignificant, consistent
209 with the idea that the spike and N gene have been under more selective constraint
210 than ORF1ab (Chaw et al., 2020). So far, no evidence has shown that vaccination or
211 other mitigation procedures causes positive selection of Omicron variants.

212

213 **Perspectives**

214 Previously we found that one to three weeks after the D value fell below -2.50, Delta
215 outbreaks emerged in India and UK, and the lambda variant outbreak in south
216 America (Yeh and Contreras, 2021b). This led to our proposal that Tajima D test
217 with a cut-off threshold value as -2.50, can predict new SARS-CoV-2 outbreak (Yeh
218 and Contreras, 2021b). Despite of the small sequence sample size (131 samples) in
219 this study, we detect a strong negative value of Tajima D in Omicron variants. This
220 finding also confirmed that the Omicron outbreak had emerged sometime before the
221 current breakout as the Tajima D test is sensitive to detect it and an efficient
222 predictor of future outbreaks. This study demonstrates that rapid genomic sequence
223 surveillance is essential, and Tajima D' tests should be included to forecast future
224 outbreaks in different geographic populations.

225

226 **Disclosure statement**

227 No potential conflict of interest was reported by the author(s).

228

229 **Author contributions**

230 All authors contributed to study concept, rationale, and initial manuscript drafts,
231 interpretation of data, and final manuscript preparation. All authors have read and
232 approved the final version of the manuscript.

233

234 **Funding**

235 No funding

236

237 **Ethical approval**

238 None declared.

239

240 **References**

241

242 Callaway, E. (2021) Heavily mutated Omicron variant puts scientists on alert.

243 Nature. <https://doi.org/10.1038/d41586-021-03552-w>

244 Chaw, S.M., Tai, J.H., Chen, S.L., Hsieh, C.H., Chang, S.Y., Yeh, S.H., Yang, W.S., Chen, P.J.,

245 Wang, H.Y. (2020) The origin and underlying driving forces of the SARS-CoV-2

246 outbreak. *J Biomed Sci.*, 27, 73.

247 Katoh, K, Standley, D.M. (2013) MAFFT multiple sequence alignment software

248 version 7: improvements in performance and usability. *Mol Biol Evol.* 30, 772–80.

249 Kuraku S, Zmasek CM, Nishimura O, Katoh K. (2013) aLeaves facilitates on-demand

250 exploration of metazoan gene family trees on MAFFT sequence alignment server

251 with enhanced interactivity. *Nucleic Acids Res.*, 41, W22-8.

252 Hahn, M.W., Rausher, M.D., Cunningham, C.W. (2002) Distinguishing between
253 selection and population expansion in an experimental lineage of bacteriophage T7.
254 *Genetics*, 161, 11-20.

255 Hughes, A.L., (2005) Evidence for Abundant Slightly Deleterious Polymorphisms in
256 Bacterial Populations. *Genetics*, 169, 533-538.

257 Hughes, A.L., Friedman, R., Rivaller, P., French, J.O., (2008) Synonymous and
258 nonsynonymous polymorphisms versus divergences in bacterial genomes. *Mol Biol*
259 *Evol.* ,25, 2199–2209.

260 Kang, L., He, G., Sharp, A.K., Wang, X., Brown, A.M., Michalak. P., Weger-Lucarelli, J.
261 (2021) A selective sweep in the Spike gene has driven SARS-CoV-2 human
262 adaptation. *Cell*, 184, 4392-4400.

263 Mugal, C.F., Wolf, J.B., Kaj, I. (2013) Why time matters: codon evolution and the
264 temporal dynamics of dN/dS. *Mol Biol Evol.*, 31,212–231.

265 Resende, P. C., Naveca, F. C., Lins, R. D. et al. (2021) The ongoing evolution of
266 variants of concern and interest of SARS-CoV-2 in Brazil revealed by convergent
267 indels in the amino (N)-terminal domain of the Spike protein. DOI:
268 10.1101/2021.03.19.21253946

269 Rozas, J., Ferrer-Mata, A., Sánchez-DelBarrio, J.C., Guirao-Rico, S., Librado, P., Ramos-
270 Onsins, S.E., Sánchez-Gracia, A. (2017) DnaSP 6: DNA Sequence Polymorphism
271 Analysis of Large Data Sets. *Mol Biol Evol.*, 34, 3299–3302.

272 Rambaut A. FigTree, version 1.4.2. San Fransisco: GitHub, Inc.; 2021. Available from:
273 <https://github.com/rambaut/figtree/releases>

- 274 Tajima, F. (1989) Statistical method for testing the neutral mutation hypothesis by
275 DNA polymorphism. *Genetics*, *123*, 585–595.
- 276 World Health Organization. (2021) Update on Omicron.
277 <https://www.who.int/news/item/28-11-2021-update-on-omicron>.
- 278 Yeh, T.Y., Contreras, G.P. (2020) Emerging viral mutants in Australia suggest RNA
279 recombination event in the SARS-CoV-2 genome. *Med J Aust* ,*213*, 44-44.e1.
- 280 Yeh, T.Y., Contreras, G.P. (2021a) Viral transmission and evolution dynamics of
281 SARS-CoV-2 in shipboard quarantine. *Bull World Health Organ.*, *99*, 486–495.
- 282 Yeh, T.Y., Contreras, G.P. (2021b) Full vaccination is imperative to suppress SARS-
283 CoV-2 Delta variant mutation frequency. *MedRxiv*,
284 <https://doi.org/10.1101/2021.08.08.21261768>
- 285 Zeng, K., Fu, Y.X., Shi, S., Wu, C.I. (2006) Statistical tests for detecting positive
286 selection by utilizing high-frequency variants. *Genetics*, *174*, 1431-1439.
- 287

288 **Figure Legend**

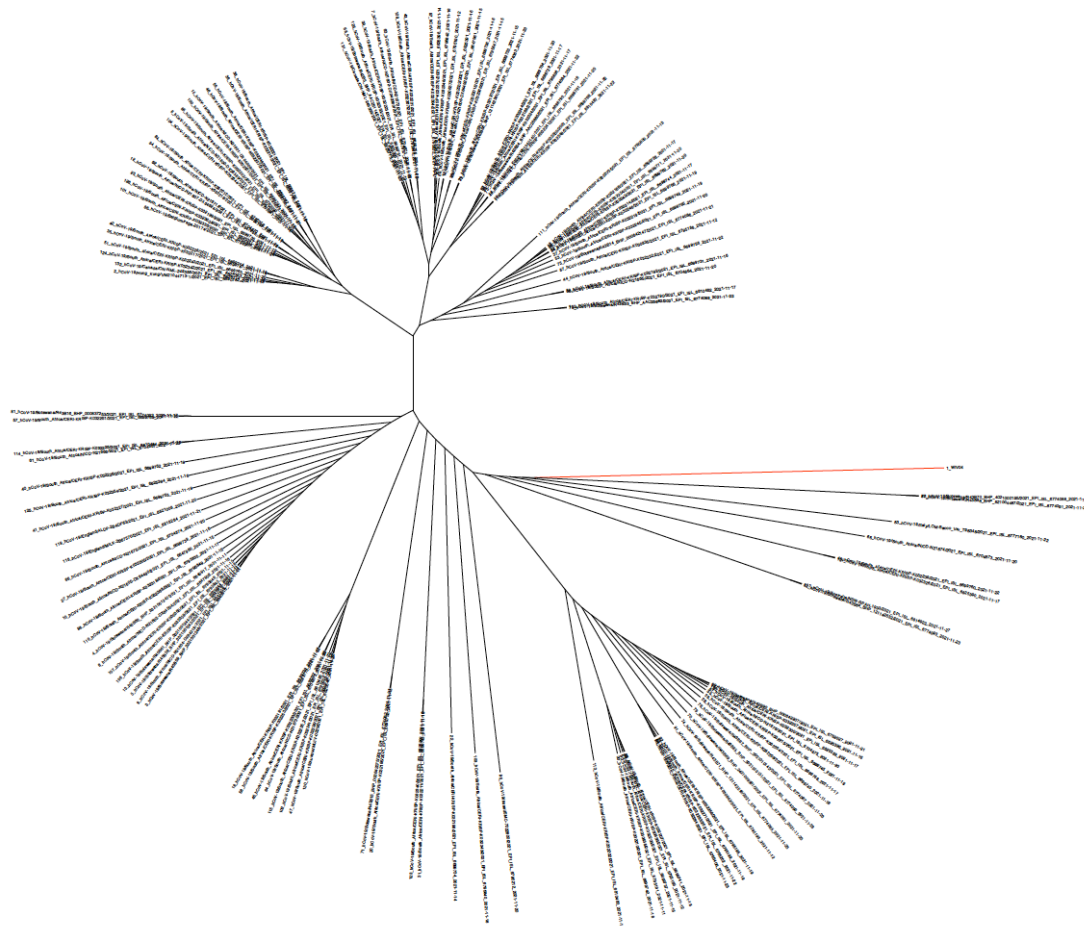
289 Figure 1. Rooted phylogenetic tree of SARS-CoV-2 genomes of Omicron variants,
290 November 9 to November 28, 2021. Alignments of viral sequences were generated
291 using MAFFT, and the phylogenetic tree was visualized using FigTree. Rooting was
292 done by introducing SARS-CoV-2 WIV04 (red line, accession number MN996528) as
293 an outgroup virus.

294 Figure 2. Alignment of Omicron variant with and without insertion at the amino acid
295 214 of the spike protein.

296

297

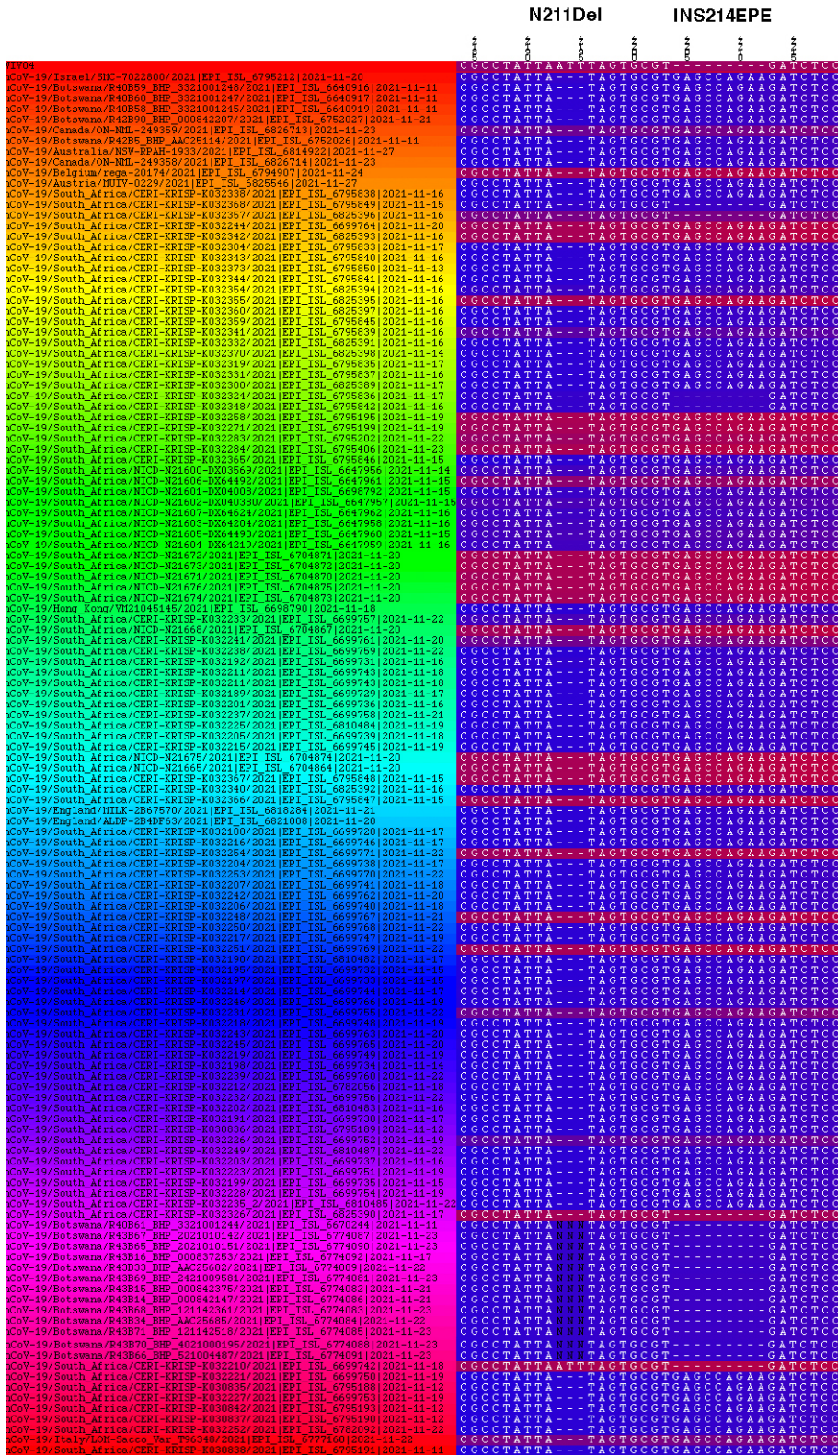
298 Figure 1.



299

300

301 Figure 2.
302



303
304

305 Table 1. Neutrality analysis of SARS-CoV-2 Omicron variants of the full genome and
 306 the viral genes. The statistical significance was estimated using 10,000 coalescent
 307 simulations in DNASP6. Not significant values ($P>0.05$) are indicated in italic. NS8 is
 308 excluded in this table because no mutations were detected. (N.D., not determined;
 309 N.A., not applicable).
 310

	Whole genome	ORF1ab	Spike	NS3	E	M	NS6	NS7a	NS7b	N
<i>DH</i>	-2.814	<i>-0.152</i>	-9.014	<i>-1.875</i>	-5.264	-7.552	<i>0.043</i>	<i>0.0588</i>	<i>0.204</i>	-3.016
Zeng's <i>E</i>	<i>0.053</i>	-2.138	5.504	<i>0.149</i>	3.272	4.603	-0.725	-0.996	-0.598	<i>0.716</i>
Tajima's <i>D</i>	-2.709	-2.617	-1.948	<i>-1.706</i>	<i>-0.735</i>	<i>-1.281</i>	-0.998	-1.342	-0.65	-2.237
<i>D</i>_{Nonsyn}	N.D.	-2.41	-1.771	-1.342	N.A.	N.A.	N.A.	N.A.	N.A.	-1.947
<i>D</i>_{syn}	N.D.	-2.46	-1.128	-1.34	N.A.	N.A.	-0.998	-1.342	-0.65	-1.57
<i>D</i> ratio (<i>D</i>_{nonsyn}/<i>D</i>_{syn})	N.D.	0.980	1.570	1.001	N.A.	N.A.	N.A.	N.A.	N.A.	1.240
ΔD (<i>D</i>_{syn}-<i>D</i>_{nonsyn})	N.D.	-0.05	0.643	0.002	N.A.	N.A.	N.A.	N.A.	N.A.	0.377