

1 **Structural Topological Analysis of Spike Proteins of**
2 **SARS-CoV-2 Variants of Concern Highlight Distinctive**
3 **Amino Acid Substitution Patterns**

4
5 Filips Peisahovics, Mohammed A. Rohaim, Muhammad Munir*

6 Division of Biomedical and Life Sciences, Lancaster University, Lancaster, Lancashire,
7 LA1 4YG, United Kingdom

8
9 ***Corresponding Author**

10
11 **Prof Muhammad Munir**

12
13 Lancaster University, Lancaster, United Kingdom.

14
15 Email: muhammad.munir@lancaster.ac.uk

33 **Abstract**

34
35 Since the onset of pandemic in 2019, SARS-CoV-2 has diverged into numerous
36 variants driven by antigenic and infectivity-oriented selection. Some variants have
37 accumulated fitness-enhancing mutations, evaded immunity and spread despite global
38 vaccination campaigns. The spike (S) glycoprotein of SARS-CoV-2 demonstrated the
39 greatest immunogenicity and amino acid substitution diversity owing to its importance
40 in the interaction with human angiotensin receptor 2 (hACE2). The S protein
41 consistently emerges as an amino acid substitution (AAS) hotspot in all six lineages,
42 however, in Omicron this enrichment is significantly higher. This study attempts to
43 design and validate a method of mapping S-protein substitution profile across variants
44 to identify the conserved and AAS regions. A substitution matrix was created based on
45 publicly available databases, and the substitution localization was illustrated on a cryo-
46 electron microscopy generated S-protein model. Our analyses indicated that the
47 diversity of N-terminal (NTD) and receptor-binding (RBD) domains exceeded that of any
48 other regions but still contained extended low substitution density regions particularly
49 considering significantly broader substitution profiles of Omicron BA.2 and BA.4/5.
50 Finally, the substitution matrix was compared to a random sample alignment of variant
51 sequences, revealing discrepancies. Therefore, it was suggested to improve matrix
52 accuracy by processing a large number of S-protein sequences using an automated
53 algorithm. Several critical immunogenic and receptor-interacting residues were
54 identified in the conserved regions within NTD and RBD. In conclusion, the structural
55 and topological analysis of S proteins of SARS-CoV-2 variants highlight distinctive
56 amino acid substitution patterns which may be foundational in predicting future variants.

57 **1. Introduction**

58 COVID-19, a disease caused by severe acute respiratory syndrome
59 coronavirus 2 (SARS-CoV-2), first appeared in the Hubei province of China in 2019.
60 On 11th March 2020, the World Health Organization (WHO) declared the disease a
61 pandemic (WHO, 2020). By 16th June 2022, the confirmed cases count surpassed 535
62 million, including 6.3 million deaths (WHO, 2022a). Furthermore, the COVID-19
63 pandemic inflicted long-term damage to several aspects of international societies,
64 including hardly quantifiable psychological impact. Yeyati and Filippini, (2021) have
65 estimated that global GDP would be 54.68% lower in 2020-2030 compared to pre-
66 pandemic trends as a consequence of educational loss, deaths and economic shrinking.
67 Governments relied on vaccination as the major countermeasure, and by June 16th,
68 2022, more than 65% of the world population (>5 billion people) had received at least
69 one dose of COVID-19 vaccine. The most commonly used vaccines around the world,
70 include Oxford-AstraZeneca, Pfizer, Moderna, Sinopharm, J&J, Sputkin-V, and Sinovac,
71 which were designed to target the spike (S) glycoprotein, which is the most immunogenic
72 protein of the virus (Das and Roy, 2021; Holder, 2022).

73 Throughout the pandemic, SARS-CoV-2 accumulated subsets of mutations
74 driven by antigenic drift and or selection favoring the virus infectivity and spread (Altmann
75 et al., 2021). WHO established Greek-letter nomenclature for variants and classified
76 them into three groups: variants under monitoring, variants of interest (VOI) and, most
77 importantly, variants of concern (VoC), which demonstrated increased transmissibility
78 and pathogenicity or decreased countermeasures efficiency (WHO, 2022b). The
79 emergence of novel variants raised concerns about vaccines efficacy, which were later
80 confirmed by a number of variants demonstrating different degrees of immune evasion
81 (Altmann et al., 2021; Jangra et al., 2021; Munir, et al., 2021; Singh et al., 2021; Wang
82 et al., 2021).

83 European Centre for Disease Prevention and Control (2022) was monitoring
84 SARS-CoV-2 variants in Europe, assessing the impact of variants' substitution portfolio
85 on severity and transmission in comparison to the previously circulating variants. Each
86 dominant variant demonstrated high transmissibility combined with immune evasion and
87 increased severity, excluding Omicron SARS-CoV-2. The variant domination patterns
88 could be judged from the representation of sequence submission dynamics to the
89 GISAID database. Alpha's transmission was surpassed by Delta variant, which was itself
90 eventually surpassed by Omicron variants (Figure 1).

91 SARS-CoV-2 is a Betacoronavirus of 65-125 nm in diameter, has positive
92 single-stranded RNA of 30-kilo base pairs genome size, encoding four structural and 15
93 accessory proteins (Jungreis et al., 2021; Astuti, 2020). S protein is a structural,
94 transmembrane glycoprotein, accommodating a homotrimer structure, each monomer-
95 1273 amino acids (141.2 kDa), its binding to the human angiotensin-converting enzyme-
96 2 receptor (hACE2) leads to viral internalisation (UniProt, 2022). The receptor-binding
97 domain (RBD) of the S protein adapted two conformations: UP – receptor accessible,
98 and DOWN – receptor inaccessible. The DOWN conformation decreased hACE2
99 recognition potential, compensated by the high affinity of RBD, and also complicated
100 antibody access (Cai et al., 2020; Shang et al., 2020). Multiple studies have reported the
101 most immunogenic and key receptor-binding residues (often overlapping) in the S protein
102 RBD, including 417, 452, 477, 484, 490, 493, 496, 498, 501 and 505, which were present
103 in several VoCs of SARS-CoV-2 (Mercurio et al., 2021; Pavlova et al., 2021; Sharma et
104 al., 2021; Watanabe et al., 2021; Yi et al., 2021; Yang et al., 2020). The S protein
105 demonstrated high diversity and mutation rates (Miao et al., 2021; Forni and Mantovani,
106 2021; Agarwal et al., 2022) such as Omicron (VOC) carried 32 mutations in S protein,
107 leading to immune evasion in vaccinated and convalescent patients (Planas et al., 2021).

108 These mutations have severely undermined vaccine efficacies and antiviral
109 therapies. The Imdevimab which targets the linear epitope (440–449 amino acid of S
110 protein), Cilgavimab and Bebtelovimab have capabilities to neutralize newest variants of
111 SARS-COV-2 including BA.2 and BA.4/BA.5 (Cao et al., 22, Ahmed et al., 2022).
112 However, antibodies such as Adintrevimab and Sotrovimab showed markedly reduced
113 neutralization against BA.4/BA.5 subvariants. Similarly, neutralization by the antibodies
114 induced by Wuhan antigen-based vaccines or through natural infections showed weaker
115 protection against Omicron subvariants particularly after four months of recovery or
116 vaccination (Cao et al., 22).

117 The aim of this study is to identify the regions of the highest immunogenicity and
118 conservation in the S protein of all major VoCs including BA.4 and BA.5 using a range
119 of *in silico* tools and models. It has been hypothesized that the functional constraint on
120 the virus divergence could result in the conservation of regions of the S protein surface
121 to preserve hACE2-interaction ability. Moreover, this study aims to apply genetic
122 analyses methods for mapping the conserved regions across all VoCs and project
123 genetic conservation in a parallel comparison fashion. The provided information is
124 fundamental to predict future evolutionary trajectories and training better vaccine and

125 therapeutic candidates.

126 **2. Methods**

127 *2.1. Sequences Acquisition*

128 Sequences pertaining to SARS-CoV-2 variants were downloaded from the
129 GISAID database in FASTA format. A query was made for each variant by inputting the
130 Pango Lineage index according to WHO and Pango Lineage nomenclature, selecting
131 a high-coverage, complete sequence without any unidentified nucleotide regions (Table
132 1).

133 *2.2. Gene Sequence Extraction and Alignment*

134 The S gene sequences was extracted from full-length sequences using
135 SnapGene Viewer software (www.snapgene.com). Each nucleotide sequence was
136 opened as a single-stranded linear DNA sequence, and then were translated to amino
137 acid sequence. Using default parameters of the SnapGene Viewer, protein-encoding
138 regions (ORF) were identified. According to Yoshimoto (2020), the S gene is located in
139 the region spanning from nucleotide 21,563 to 25,384. Both nucleotide and amino acid
140 sequences were extracted and saved as FASTA format in individual datasets. Next, the
141 translated amino acid sequences and a Wuhan reference S protein amino acid
142 sequence were aligned using MEGAX software by the MUSCLE method (with
143 parameters including Gap open: -2.9; Gap extend: 0; Hydrophobicity multiplier: 1.2;
144 Max iterations: 16; Cluster method: UPGMA; Lambda: 24) (Edgar, 2004) and the
145 alignment was exported in FASTA format.

146

147 *2.3. Analysis of Amino Acid Substitutions in the S protein*

148 A table of amino acid substitutions (data not presented) was made for each
149 variant according to the CoVariants Online Database (Hodcroft, 2022), excluding Zeta
150 and Theta, that were acquired directly from GISAID from the earliest accession in the
151 Pango lineage-filtered tree (Khare et al., 2021). The substitution recurrence across
152 variants (cumulative substitution count) was calculated per each identified site and
153 summed by the S protein regions according to the GenBank designation (Accession ID:
154 NC_0455122). The relative frequency was calculated as the ratio of cumulative
155 substitution count in sites of a region against the total number of substitutions in the
156 protein. These results were organized in a substitution matrix format.

157 The substitutions present in the alignment were compared against the
158 substitution matrix to identify discrepancies. To construct substitution sequence
159 WebLogo representation, a sequence set of varying residues (according to the

160 substitution matrix) was constructed in FASTA format and imported into the WebLogo
161 online tool (Crooks et al., 2004). Finally, the alignment was uploaded to the Phylogeny.fr
162 online tool (“One Click” Mode) to construct a maximum-likelihood phylogenetic tree
163 (Dereeper et al., 2008).

164

165 *2.4. Spike protein Substitutions 3D*

166 A cryo-electron microscopy 3D structure model of the Wuhan’s S protein with
167 RBD in the UP conformation was downloaded from RCSB Protein Data Bank
168 (accession ID: 6SVB) (Wrapp et al., 2020). Each variant was designated with a random,
169 visually-distinguishable color to help in the comprehension of the models. All the
170 substitution sites were highlighted at one subunit of the Wuhan S protein using UCSF
171 Chimera software. Despite every substitution was being highlighted with a contrasting
172 color, some of the substitutions were not visible from the chosen perspective, therefore,
173 not indicated.

174 Next, 3D models were generated for each variant, highlighting the substitutions
175 localisation on every subunit. Substitution sites were highlighted according to the
176 substitution matrix; if the substitution residue was not resolved in the model, the two
177 flanking residues were highlighted, unless one of those was above ten residues apart,
178 then only the closer residue was highlighted.

179

180 **3. Results**

181 *3.1. Substitution Localisation*

182 Our results revealed that the S protein fusion peptide (798-806), internal fusion
183 peptide (816-833), transmembrane domain (1213-1236) and cytoplasmic domain
184 (1237-1273) regions did not have any reported substitutions in the sample alignment.
185 The majority of the reported substitutions were localized in the N-terminal half-part of
186 the S protein, especially in the N-terminal domain (NTD) and RBD (Figure 2A). The C-
187 terminal half-part contained the minor part of the substitutions that were localised in four
188 function-unidentified intermediate regions (IR), heptad repeats (HR) I and II, and the
189 cleavage site (CS) (Figure 2B). In addition, the level of conservation varied across
190 regions (Figure 2).

191 Although, NTD has the highest number of substitutions, the consensus sequence
192 remained considerably conserved compared to other regions. Meanwhile, RBD and CS
193 showed the highest variability; particularly, RBD has highly variable sites 417, 452, 484
194 and 501, while CS has sites 681 and 701. Substitution sites 142, 253, 339, 371, 655

195 and 1176 also demonstrated outstanding variability, yet to a lower extent compared to
196 the aforementioned (Figure 3).

197

198 3.2. *Phylogenetic Analysis*

199 Maximum likelihood tree for the S protein amino acid sequences as well as
200 phylogenetic tree based on full length sequence of the SARS-COV-2 revealed that
201 Omicron BA.1, BA.2, BA.4 and BA.5 subvariants were significantly deviated from all
202 other variants (Figure 4, A, B). Notably, the variant clustering did not correspond to the
203 chronological outbreak order (WHO, 2020b), while the Wuhan clustering closely with
204 Lambda, Delta, Kappa and Epsilon. Additionally, all variants including Omicron sub-
205 variants formed a successive, single branching pattern; every earlier branch diverged
206 only once, into a single later branch (Figure 4, A and B).

207

208 3.3. *Cumulative Substitutions*

209 The proportion of cumulative substitutions prevailed in NTD and RBD – 33.88%
210 and 35.95% respectively. Moreover, out of total 93 overall substitution sites, RBD and
211 NTD contained 66 substitutions. The RBD showed a greater number of substitutions
212 than NTD (87 versus 82 subs) and more substitutions per residue (0.39 VS 0.28) due
213 to the shorter length of the region (223 versus 292 amino acids). Out of 93 unique
214 substitution sites, most of these substitutions were present in one to few variants where
215 47 sites were unique to one variant, 34 were present in two to four, and the remaining
216 12 sites – in five or more variants. Several sites have an outstanding manner of
217 substitutions: site 142 and 144 – in 5 variants, site 417 – in 6, sites 501 and 681 – in 9
218 variants, site 484 – in 12 variants, and all 16 variants have an Asp-to-Gly substitution
219 at site 614 (Table 2).

220

221 3.4. *Substitution Patterns*

222 3.4.1. Substitutions in NTD

223 Deletions constitute around 50% of the total amino acid substitutions in the
224 NTD where some of these deletions were observed in several SARS-CoV-2 variants.
225 Variants including Alpha, Eta and Omicron BA.1, BA.4 and BA.5 have H69del and
226 V70del deletions. One specific deletion (Y144del) also appeared in Alpha, Eta,
227 Omicron BA.1 as well as within a deleted region 141-144 in Theta. Similarly, region
228 from 142-145 amino acids was substituted in Omicron BA.1 while Alpha and Theta
229 variants shared a pair of identical deletions; L241del and A243del. Lambda variant

230 showed the highest number deletions as a distinct, where seven amino acids long
231 sub-region, spanning from site 246 to 252 inclusively, along with substitution D253N,
232 followed by Theta and Omicron BA.1 variants with six deletions, and Omicron BA.4
233 and BA.5 with five deletions (Table 2). The most common substitution in NTD, apart
234 from deletions, was T95I, appearing in Iota, Mu and Omicron BA.1, while Omicron
235 subvariants BA.2, BA.4 and BA.5 also shared T19I, A27S and G142D mutations.
236 Finally, Lambda and Omicron BA.1 showed the highest overall number of
237 substitutions in NTD (ten substitutions), followed by Omicron BA.4 and BA.5 (nine
238 substitutions). In contrast, Zeta carried no mutations in the NTD while Kappa has only
239 one, and Epsilon and Iota have two substitutions (Table 2).

240

241 3.4.2. Substitutions in RBD

242 Omicron subvariants carried 15-17 substitutions in the RBD, significantly
243 outperforming all other variants, which accumulated only one to three substitutions in
244 the RBD of S-protein. Interestingly, just one substitution (F486V) of total 17 was unique
245 to the most-recent Omicron subvariants BA.4 and BA.5, which shared identical
246 substitution portfolio in S protein. Apart from Omicron-specific substitutions, most of the
247 substitutions in RBD were present in multiple variants, and R364K and F490S
248 substitutions appeared in Mu and Lambda variants.

249 Residue 484 is the highest diversified among all the reported substitution sites
250 within the S-protein. Alpha, Delta, Epsilon and Lambda variants retained the original
251 Wuhan's Glu at this site, while other variants, except for both Omicrons and Kappa were
252 substituted by Lys – the most common form across variants. Omicrons have Ala at this
253 site, but Kappa has a unique Gln. Six variants that have accumulated N501Y
254 substitution included Alpha, Beta, Gamma, Theta, Mu and Omicron subvariants,
255 making it the second most common substitution in the RBD, followed by L452R that is
256 in six variants. Notably, substitutions L452R and E484(K/A) did not overlap (Table 2).

257

258 3.4.3. Substitutions in Other Regions

259 The remaining regions, including SP, IR 1-4, S1/S2 CS and HR I and II, carried
260 additional unique 27 substitution sites out of a total of 93 substitutions. Disrespecting
261 Omicron-specific substitutions, only six were present in two or more variants, including
262 (1) D614G that as present in every variant, (2) H655Y, (3) P681 (H/R) and (4) A701V
263 in S1/S2 CS, (5) D950N in HR I and (6) V1176F in HR II.

264 Residue 681 showed the second-highest diversified site after residue 484 in the

265 RBD. Variants Alpha, Theta, Mu and all Omicron subvariants carried a Pro-to-His
266 substitution at this site, while Delta and Kappa have a Pro-to-Arg one. Omicron
267 subvariants demonstrated the greatest substitution numbers in the remaining regions –
268 8 in BA.2, BA.4 and BA.5, and 10 in BA.1. Variants such as Alpha and Theta revealed
269 six substitutions in these regions while Gamma contains four, and others - one to three
270 (Table 2).

271

272 3.4.4. Conserved Regions

273 Regions within the S-protein were only partially diversified while NTD spanned
274 from residue 13 to 304. There were numerous subregions of low substitution density,
275 including subregions 27-67 that contain a single cumulative substitution (Q52R), 95-
276 138, 158-211 (R190S) and 215-241. Next, inter-region 265-339 located between NTD
277 and RBD (spanning 319-541) has no substitutions. Furthermore, from site 265 to 371,
278 there were only two substituted residues; G339D and R346K. Going downstream of
279 the RBD, several subregions were of a low substitution density: 376-405; 417-440;
280 and 505-547. The longest substitution-less subregion of the NTD- RBD region was
281 the inter-region 265-339 (74 amino acids long), followed by 95-138 (43), 505-547 (42)
282 and 339-371 (32). Additionally, subregions 158-211 (58) and 27-67 (40) contain only
283 one substitution (Table 2).

284

285 3.4.5. Sample Alignment and Matrix Discrepancies

286 Most of the reported substitutions were present in the sample alignment with
287 four exceptions, and additional substitutions were discovered in six cases. Sequences
288 from Beta, Iota, Kappa and Lambda contain two additional substitutions each, while
289 Mu had three, and Delta – four; also, Delta and Theta sequences were missing two
290 reported substitutions each, Mu missed one, and Omicron BA.1 was missing eight
291 reported substitutions in the NTD. In contrast, sequences of Alpha, Gamma, Epsilon,
292 Zeta, Eta and Omicron BA.2, BA.4 and BA.5 corresponded to the reported
293 substitutions profile perfectly (Table 3).

294

295 3.4.6. Substitutions Pattern Visualization

296 The localization of substitutions on the surface of S-protein subunits is shown in
297 Figure 4. Overall, 23 out of 93 substitution sites were not visible on the surface of the
298 S-protein from three selected perspectives. The top-centre region (Figure 7) of the
299 protein surface contains a dense substitution localization that accommodated most of

300 the substitution sites in the RBD, except for 405, 408 and 446. However, RBD sub-
301 regions 471-487 and 371-376 were unresolved in the cryo-EM structure, hence, the
302 surface protrusion of the contained six substitution residues could not be visualized
303 precisely.

304 The pyramidal, outer-protruding structures at the topsides of the protein
305 demonstrated substitution sites in the NTD, missing sites up to site 26 due to unresolved
306 N-terminal gap in the cryo-EM structure as well as substitution site 265. The NTD
307 subregions including 66-80, 140-158, and 245-261 flanked 22 substitution sites, whose
308 positioning in the cryo-EM model was unresolved. The remaining 15 surface-visible
309 sites were distributed across other S-protein region that was downstream in primary
310 sequence, and ten other sites in these regions were not visible on the surface. Sites
311 677, 679 and 681 were unresolved and, therefore, indicated by a flanking subregion
312 672-687.

313 In comparison to substitution clustering in the RBD and NTD, the remaining
314 substitutions were distributed individually. Most of these substitutions were localised on
315 the top surface (Figure 7) were common across variants with a few exceptions: R346K
316 in Mu, T376A in Omicron BA.2, K417T in Gamma, E484Q in Kappa, F490K in Lambda
317 and G496S in Omicron BA.1, BA.4 and BA.5. Remarkably, the highest substitutions
318 density across variants was localized in the centre of the top surface, where the RBD
319 (Figure 7 and 8). Omicron subvariants have the highest substitution density on the top
320 surface. Uniquely, variant-specific mutations on the S-protein sides mainly appeared in
321 the NTD, particularly in unresolved subregions spanning from amino acid 140-157 and
322 245-261. Overall, the substitution localization density decreased from the RBD on the
323 top surface down to the cytosolic domain (Figure 7 and 9).

324

325 **4. Discussion**

326 Early investigations of SARS-CoV-2 genomes predicted an evolutionary rate of
327 roughly 0.001 subs site⁻¹ year⁻¹ (two to three mutations per month) (Duchene et al.,
328 2020); however, there is significant divergence from this pace across the phylogeny,
329 with certain outlier lineages, particularly VOC, acquiring multiple mutations at a
330 considerably faster rate. The mutations analysis from virus genome data is important
331 for basic virology (Houldcroft et al. 2017), as it identifies evolutionary signals associated

332 with mutations prior to experimental and real-world data on clinical outcomes or vaccine
333 effectiveness, and it documents and tracks changes that may affect therapeutic
334 effectiveness. Therefore, it is imperative to assess the tendencies and trends in the
335 topological and structural differences of major variants of concerns to predict future
336 evolutionary trajectories (Tariq et al., 2022).

337 Currently, about 12 million genome sequences are accessible through the
338 GISAID Initiative, allowing for real-time monitoring of the epidemic (Shu and McCauley
339 2017; Meredith et al. 2020). Since the cumulative substitution count was based on the
340 number of substitution recurrences in WHO-named variants. Theoretically, it indicates
341 the importance of the S protein in terms of phenotypical diversification, gaining
342 enhanced transmissibility, pathogenicity, evading immunity or adapting for a particular
343 epidemiological niche (Wright et al., 2022). Hence, a high cumulative substitution count
344 of a region would suggest that substitution accumulation in this region contributed to
345 the emergence of SARS-CoV-2 variants with the aforesaid qualities more substantially
346 than substitution accumulation in regions of a low cumulative substitution count.

347 Both NTD and RBD of the S protein demonstrated the highest cumulative
348 substitution count – 82 (33.88%) and 87 (35.95%), respectively, and contained the
349 majority of all substitution sites – 66 out of 93. Multiple sites augmented the diversity
350 of RBD by varying substitutional outcomes, such as Glu-to-Lys, Glu-to-Gln and Glu-to-
351 Ala substitutions emerged at site 484 in various variant combinations. Besides,
352 deletions prevailed in the NTD, making 43 out of 82 substitutions in total. Finally, the
353 visualization of substitution on cryo-EM models suggested that the substitution density
354 prevailed in the RBD and NTD, the distal-most domains of the S-protein, which were
355 reported to be targeted by antibodies (Dejnirattisai et al., 2021; McCallum et al., 2021;
356 Liu et al., 2020; [Ahmed et al., 2022](#)). Thus, RBD and NTD, to a lesser extent, could be
357 considered as the critical S protein region regarding viral adaptation, immune response,
358 and treatment design.

359 Omicrons shared substitutions that were associated with enhanced hACE2
360 affinity and immune evasion as in Alpha variant - N501Y and P681H (Khan et al., 2021;
361 Luan et al., 2021). While Omicrons had one shared substitution in RBD of unclear
362 importance with Delta, thought to be linked to viral fitness - T478K (Di Giacomo et al.,

363 2021; Jhun et al., 2021). Despite substitution E484K being reported to have a role in
364 evading immunity (Wu et al., 2022; Jangra et al., 2021), neither Alpha nor Delta
365 contained it, while it reported in seven other variants. Omicrons have an E484A
366 substitution, which similarly to E484K, removed the carboxyl's negative charge.
367 Perhaps, the removal of negative charge at site 484 resulted in the immune evasion by
368 decreasing antibody recognition ability, but the introduction of Lys positive charge
369 created more epitope potential than the introduction of neutral Ala residue as in
370 Omicrons (Wu et al., 2022; Jangra et al., 2021; Altaf et al., 2022). Interestingly, recently
371 emerged Omicron variants consists of identical RBD sequences compared to BA.2 with
372 the L452R/F486V mutations in BA.4 and BA.5. These mutations provided a
373 transmission advantage and antibodies escape characteristics to BA.4 and BA.5 over
374 BA.2. This aligned with the recent research highlighting those individuals who had
375 recovered from SARS infection displayed a systematic reduction in neutralization
376 activity against BA.4 and BA.5 variants (Cao et al., 2022; Agarwal et al., 2022). Overall,
377 Omicron subvariants' substitution profile was explicitly distinct from others, clearly
378 illustrated by the phylogenetic analyses. The unprecedented number of substitutions in
379 the RBD could explain the enhanced transmissibility and immune evasion ability of
380 Omicron variants (ECDC, 2022; Planas et al., 2021).

381 Strong dependence on the host's protein machinery set functional constraints on
382 the adaptive capability of RNA viruses, leaving a genetic vulnerability as RNA viruses
383 could keep substituting only a fraction of their genome to evade immunity before being
384 forced out of the niche (Simmonds et al., 2019; Holmes, 2003). This functional
385 constraint might press on the SARS-CoV-2 to conserve particular on critical regions in
386 the S protein throughout variants to be able to bind with hACE2 – these regions were
387 arguably the low substitution density subregions of NTD and RBD.

388 Several studies have reported immunogenic and receptor-binding key residues
389 in the RBD. Six of them were not substituted in any of the observed variants and situated
390 in the extended, low substitution density subregions of RBD: 403, 418, 421, 426, 439
391 and 506 (Pavlova et al., 2021; Watanabe et al., 2021; Yi et al., 2021; Yang et al., 2020).
392 Additionally, Mercurio et al. (2021) and Sharma et al. (2021) have reported other critical
393 residues located in shorter low substitution density regions of the RBD. Despite eliciting

394 strong immune response, these residues were not substituted, which might indicate that
395 these residues and corresponding subregions were under functional constraint,
396 therefore, treatment targeted them would potentially be less variant-biased. Even highly
397 mutated Omicron subvariants contained low substitution density subregions in the RBD,
398 where the six key residues remained unsubstituted.

399 Our analyses revealed that individual virions attributed to a variant could contain
400 additional substitutions and lack those reported to be variant-specific, so the variant
401 designation might only represent a particular fraction of viruses in the lineage. Besides,
402 previous study estimated the variant doubling period at 71.67 (\pm 0.06 SE), one novel
403 variant per 600,000 infections, which signified the divergence potential of SARS-CoV-
404 2 and, hence, the necessity for universal treatment (Duarte et al., 2021).

405 Results of substitution matrix and cumulative substitution count support the
406 importance of NTD and RBD in the antigenic drift of SARS-CoV-2. Additionally, the
407 substitution matrix clearly illustrated the differences between substitution profiles
408 across variants, especially of Omicron subvariants, which could be studied further
409 regarding adaptation patterns, biochemical and epidemiological effects. Overall, this
410 study has demonstrated a method of genetic analysis that would hypothetically aid in
411 revealing sites and regions in the S protein of high immunogenicity and conservation
412 if improved in terms of scope and accuracy.

413 We offered in this study an algorithm that would compute a substitution matrix
414 with cumulative substitution count per residue based on the large alignment of all
415 SARS-CoV-2 sequences. In addition, the hypervariable residues could be mapped
416 using the cryo-EM model to check the antibody binding sites and corrected for glycan
417 shielding that can cover 40% of the S-protein (Grant et al., 2020). However, an
418 enormous computational capacity would be required to analyze such query, as only
419 for the alignment stage, the best US supercomputer operated a full week to align only
420 17,000 virus genomes (Garvin et al., 2020), and the sequence number is the primary
421 accuracy-limiting factor. Therefore, to facilitate such calculation, the amino acid
422 sequences can be trimmed down to include NTD and RBD due to their immunogenicity
423 and receptor-interaction importance.

424

425 **5. Conclusions and limitations**

426 Tracking the virus evolution play an important role in providing clear and
427 accessible information to those who are tackling the pandemic, including through public
428 health actions and the development of vaccines and therapeutics. Although amino acid
429 sequence analyses alone are insufficient to determine the functional effect of a single
430 mutation on SARS-CoV-2 fitness, computational analysis of existing SARS-CoV-2
431 mutations provided substantial information on possible phenotypic changes and
432 projected mutations may confer on variants.

433

434 **CRedit authorship contribution statement**

435 Filips Peisahovics, Mohammed A. Rohaim and Muhammad Munir: Conceptualization,
436 Methodology, Investigation, Writing – original draft, Visualization. Muhammad Munir:
437 Resources, Data curation. Filips Peisahovics and Mohammed A. Rohaim, Writing –
438 review & editing. Muhammad Munir Supervision and Funding Acquisition. All authors
439 read and confirmed the submission of the paper.

440

441 **Data Availability**

442 Data will be made available on request.

443

444 **Acknowledgement or Funding**

445

446 *Declarations of interest*

447 None

448

449 **References**

450

451 Agarwal, D, Zafar, I, Ahmad, S.U, et al., 2022. Structural, genomic information and
452 computational analysis of emerging coronavirus (SARS-CoV-2). Bull Natl Res Cent.
453 46(1): 170.

454 Ahmad, S.U, Kiani, B.H, Abrar, M. et al., 2022. A comprehensive genomic study,
455 mutation screening, phylogenetic and statistical analysis of SARS-CoV-2 and its variant
456 omicron among different countries. J Infect Public Health. 15(8):878-891.

457 Altmann, D.M., et al., 2021. Immunity to SARS-CoV-2 variants of concern. *Science*,

458 371(6534): p. 1103-1104.

459 Astuti, I., 2020. Severe Acute Respiratory Syndrome Coronavirus 2 (SARS- CoV-2):
460 An overview of viral structure and host response. *Diabetes & Metabolic Syndrome:
461 Clinical Research & Reviews*, 14(4): 407–412

462 Cai, Y., et al., 2020. Distinct conformational states of SARS- CoV-2 spike protein.
463 *Science*, 369(6511): 1586-1592

464 Cao, Y., Yisimayi, A., Jian, F. et al. BA.2.12.1, BA.4 and BA.5 escape antibodies
465 elicited by Omicron infection. *Nature* (2022). [https://doi.org/10.1038/s41586-022-
466 04980-y](https://doi.org/10.1038/s41586-022-04980-y)

467 Crooks, G.E., et al., 2004. WebLogo: a sequence logo generator. *Genome Research*,
468 14(6): 1188-1190

469 Altaf M., et al., 2022. Wildlife as a Source of SARS-CoV-2 Evolution- A Review.
470 *Pakistan journal of Zoology*. 54 (4): 1899-1904

471 Das, J.K. and Roy, S., 2021. A study on non-synonymous mutational patterns in
472 structural proteins of SARS-CoV-2. *Genome*, 99(999): 1-14

473 Dejnirattisai, W., et al., 2021. The antigenic anatomy of SARS-CoV-2 receptor binding
474 domain. *Cell*, 184(8): 2183-2200

475 Dereeper, A., et al., 2008. Phylogeny. fr: robust phylogenetic analysis for the non-
476 specialist. *Nucleic acids research*, 36(2): 465-469

477 Di Giacomo, S., et al., 2021. Preliminary report on severe acute respiratory syndrome
478 coronavirus 2 (SARS- CoV-2) Spike mutation T478K. *Journal of medical virology*, 93(9):
479 5638-5643

480 Duarte, C.M., et al. 2022. Rapid evolution of SARS-CoV-2 challenges human defenses.
481 *Scientific Reports*, 12(1); 6457.

482 Edgar, R.C., 2004. MUSCLE: multiple sequence alignment with high accuracy and high
483 throughput. *Nucleic acids research*, 32(5): 1792-1797.

484 European Centre for Disease Prevention and Control (ECDC), 2022. SARS-CoV-2
485 variants of concern as of 24 February 2022. Available at: [https://www.ecdc.europa.
486 eu/en/covid-19/variants-concern](https://www.ecdc.europa.eu/en/covid-19/variants-concern) [Accessed on 1st March 2022].

487 Forni, G. and Mantovani, A., 2021. COVID-19 vaccines: where we stand and

488 challenges ahead. *Cell Death & Differentiation*, 28(2): 626-639.

489 Garvin, M.R., et al., 2020. A mechanistic model and therapeutic interventions for
490 COVID-19 involving a RAS-mediated bradykinin storm. *eLife*, 9: e59177.

491 Grant, O.C., et al., 2020. Analysis of the SARS-CoV-2 spike protein glycan shield
492 reveals implications for immune recognition. *Scientific reports*, 10(1): 1-11.

493 Jangra, S., et al., 2021. The E484K mutation in the SARS-CoV-2 spike protein reduces
494 but does not abolish neutralizing activity of human convalescent and post-vaccination
495 sera. *MedRxiv*.01.26.21250543. doi: 10.1101/2021.01.26.21250543.

496 Jhun, H., et al., 2021. SARS-CoV-2 Delta (B.1.617.2) variant: a unique T478K mutation
497 in receptor binding motif (RBM) of spike gene. *Immune Network*, 21(5): 32.

498 Jungreis, I., et al., 2021. Conflicting and ambiguous names of overlapping ORFs in the
499 SARS-CoV-2 genome: A homology-based resolution. *Virology*, 558: 145-151

500 Hodcroft, E.B., 2021. CoVariants: SARS-CoV-2 Mutations and Variants of Interest.
501 Available at: <https://covariants.org/> [Accessed on 14th June 2022].

502 Holder, J., 2022. Tracking coronavirus vaccinations around the world. Available at:
503 <https://www.nytimes.com/interactive/2021/world/covid-vaccinations-tracker.html>
504 [Accessed on 16th June 2022].

505 Holmes, E.C., 2003. Error thresholds and the constraints to RNA virus evolution. *Trends*
506 *in microbiology*, 11(12): 543-546.

507 Khan, A., et al., 2021. Higher infectivity of the SARS-CoV-2 new variants is associated
508 with K417N/T, E484K, and N501Y mutants: an insight from structural data. *Journal of*
509 *cellular physiology*, 236(10): 7045-7057.

510 Khare, S., et al., 2021. GISAID's Role in Pandemic Response. *China CDC Weekly*,
511 3(49): 1049.

512 Liu, L., et al., 2020. Potent neutralizing antibodies against multiple epitopes on SARS-
513 CoV-2 spike. *Nature*, 584(7821): 450-456.

514 Luan, B., et al., 2021. Enhanced binding of the N501Y- mutated SARS-CoV-2 spike
515 protein to the human ACE2 receptor: insights from molecular dynamics simulations.

516 *FEBS letters*, 595(10): 1454-1461.

517 McCallum, M., et al., 2021. N-terminal domain antigenic mapping reveals a site of
518 vulnerability for SARS-CoV-2. *Cell*, 184(9): 2332-2347.

519 Mercurio, I., et al., 2021. Protein structure analysis of the interactions between SARS-
520 CoV-2 spike protein and the human ACE2 receptor: from conformational changes to
521 novel neutralizing antibodies. *Cellular and Molecular Life Sciences*, 78(4): 1501-1522.

522 Miao, M., et al., 2021. Genetic Diversity of SARS-CoV-2 over a One-Year Period of the
523 COVID-19 Pandemic: A Global Perspective. *Biomedicines*, 9(4): 412.

524 Nikolaidis, M.; et al., 2022. Comparative Analysis of SARS-CoV-2 Variants of Concern,
525 Including Omicron, Highlights Their Common and Distinctive Amino Acid Substitution
526 Patterns, Especially at the Spike ORF. *Viruses*, 14, 707.

527 Pavlova, A., et al., 2021. Machine Learning Reveals the Critical Interactions for SARS-
528 CoV-2 Spike Protein Binding to ACE2. *The journal of physical chemistry letters*, 12(23):
529 5494-5502.

530 Munir M., A. et al., 2021. Facts and Figures On Covid-19 Pandemic Outbreak.
531 *Pakistan journal of Zoology*, 53 (3): 801-1200

532 Tariq, M. H et al., 2022. In silico Screening of Bioactive Phytochemicals against Spike
533 Protein of COVID-19. *Pakistan journal of zoology*, 54 (1): 433-438

534 Planas, D., et al., 2021. Considerable escape of SARS-CoV-2 Omicron to antibody
535 neutralization. Available at: [https:// www.nature.com/articles/s41586-021-04389-z](https://www.nature.com/articles/s41586-021-04389-z)
536 [Accessed on 10th February 2022].

537 Shang, J., et al., 2020. Cell entry mechanisms of SARS-CoV-2. *Proceedings of the*
538 *National Academy of Sciences*, 117(21): 11727-11734.

539 Sharma, D., et al., 2021. Elucidating important structural features for the binding affinity of
540 spike-SARS-CoV-2 neutralizing antibody complexes. *Proteins: Structure, Function,*
541 *and Bioinformatics*, 90(3): 824- 834.

542 Simmonds, P., et al., 2019. Prisoners of war - host adaptation and its constraints on
543 virus evolution. *Nature Reviews Microbiology*, 17(5): 321-328.

544 Singh, J., et al., 2021. Evolutionary trajectory of SARS-CoV-2 and emerging variants.
545 *Virology journal*, 18(1): 1-21.

546 Universal Protein Database, 2022. P0DTC2 · SPIKE_SARS2. Available at:
547 <https://covid-19.uniprot.org/uniprotkb/P0DTC2> [Accessed on 10th February 2022].

548 Wang, L., et al., 2021. Ultrapotent antibodies against diverse and highly transmissible
549 SARS-CoV-2 variants. *Science*, 373(6556): 818-823.

550 Watanabe, K., et al., 2021. Intermolecular Interaction Analyses on SARS-CoV-2 Spike
551 Protein Receptor Binding Domain and Human Angiotensin-Converting Enzyme 2
552 Receptor-Blocking Antibody/Peptide Using Fragment Molecular Orbital Calculation. *The*
553 *journal of physical chemistry letters*, 12(16): 4059-4066.

554 Wise, M.J., 2020. SARS, SARS-Cov-2 and MERS are the Same Viral Species, Clades
555 within Bat Beta-Coronaviruses, *OSF Preprints*.

556 World Health Organisation, 2022a. WHO Coronavirus (COVID-19) Dashboard. Available
557 at: <https://covid19.who.int/> [Accessed on 16th June 2022].

558 World Health Organisation, 2022b. Tracking SARS-CoV-2 variants. Available at:
559 <https://www.who.int/en/activities/tracking-SARS-CoV-2-variants/> [Accessed on 11th
560 February 2022].

561 World Health Organisation, 2020. WHO Director-General's opening remarks at the
562 media briefing on COVID-19 - 11 March 2020. Available at: [https://www.who.int/director-
563 general/speeches/detail/who-director-general-s-opening-remarks-at-the-media-
564 briefing-on-covid-1911-march-2020](https://www.who.int/director-general/speeches/detail/who-director-general-s-opening-remarks-at-the-media-briefing-on-covid-1911-march-2020) [Accessed on 9th February 2022]

565 Wrapp, D., et al., 2020. Cryo-EM structure of the 2019-nCoV spike in the prefusion
566 conformation. *Science*, 367(6483): 1260-1263.

567 Wright, D.W., et al., 2022. Tracking SARS-CoV-2 mutations and variants through
568 the COG-UK-Mutation Explorer. *Virus evolution*, 8(1), veac023.
569 <https://doi.org/10.1093/ve/veac02>

570 Wu, L., et al., 2022. Exploring the immune evasion of SARS-CoV-2 variant harboring
571 E484K by molecular dynamics simulations. *Briefings in bioinformatics*, 23(1): 383.

572 Yang, Y., et al., 2021. Key residues of the receptor binding domain in the spike protein
573 of SARS- CoV-2 mediating the interactions with ACE2: a molecular dynamics study.
574 *Nanoscale*, 13(20): 9364-9370.

575 Yeyati, E.L. and Filippini, F., 2021. Social and economic impact of COVID-19. Available
576 at: <https://www.brookings.edu/research/social-and-economic-impact-of-covid-19/>
577 [Accessed on 9th February 2022].

578 Yi, C., et al., 2020. Key residues of the receptor binding motif in the spike protein of
579 SARS-CoV-2 that interact with ACE2 and neutralizing antibodies. *Cellular & molecular*
580 *immunology*, 17(6): 621-630.

581 Yoshimoto, F.K., et al., 2020. The proteins of severe acute respiratory syndrome
582 coronavirus-2 (SARS-CoV-2 or n-COV19), the cause of COVID-19. *The protein journal*,
583 39(3): 198-216.

584 **Figure legends:**

585 **Figure 1:** SARS-CoV-2 variants distribution in sequences submitted to GISAID
586 database.

587 **Figure 2:** Localization of cumulative substitutions in regions of SARS-CoV-2 spike
588 protein. (A) Protein regions correspond to GenBank designation (Accession ID:
589 NC_045512.2). Cumulative substitutions were calculated as the sum of recurrence
590 events in all the variants. The relative frequency (presented in brackets) represents the
591 ratio of the cumulative substitution frequency in the sites of a region against the
592 cumulative number of substitutions in the protein. (B) Schematic domain structure of S
593 protein. Different domains including signal peptide (SP), N-terminal domain (NTD)
594 receptor binding domain (RBD), spike protein subunit 1 (S1), spike protein subunit 2
595 (S2), fusion peptide (FP), heptad repeat 1 domain (HR1), heptad repeat 2 domain
596 (HR2), and transmembrane domain (TM) are shown.

597 **Figure 3:** SARS-CoV-2 S protein substitution sequence WebLogo representation.
598 Protein regions are indicated above the residues according to the GenBank designation
599 (Accession ID: NC_045512.2). Abbreviations: CS - Cleavage Site; HR - Heptad Repeat;
600 IR - Intermediate Region; SP - Signaling Peptide. (Adapted from: Crooks et al., 2004).

601 **Figure 4:** Phylogenetic analysis of SARS-CoV-2 variants based on Spike protein
602 amino-acid sequence (A) and full-length sequences of the SARS-COV-2 (B).
603 Phylogeny.fr online application was used to construct a maximum-likelihood tree.
604 Sample sequences were obtained from the GISAID database, accession IDs presented
605 in the branch names.

606 **Figure 5:** SARS-CoV-2 variants colour designation. The colour choice was completely
607 random. This figure was designed to aid comprehension of the following figures.

608 **Figure 6:** Localisation of amino acid substitution sites on the surface of SARS-CoV-2
609 spike protein. Presented cryo-electron microscopy structure was obtained from Protein
610 Data Bank (Accession ID: 6VSB) and processed in UCSF Chimera software. The
611 receptor-binding domain in this model adapted an UP conformation. Three subunits are
612 colored in shades of red, and the most saturated one is the subunit of interest, on which
613 the substitution localisation was highlighted according to those reported at least in one
614 Greek letter-designated variant.

615 **Figure 7:** Comparison of amino acid substitution localisation in spike proteins of SARS-
616 CoV-2 variants. Presented cryo-electron microscopy structure was obtained from
617 Protein Data Bank (accession ID: 6VSB) and processed in UCSF Chimera software.
618 The receptor-binding domain in this model adapted an UP conformation. The
619 localisation of reported substitution residues was highlighted on all three subunits (see
620 Section 2.3). Unique, variant-specific substitutions were indicated. Panel A
621 demonstrates top perspective, Panel B - side-front, and Panel C - another side.

622 **Figure 8:** Comparison of amino acid substitution localisation in spike proteins of SARS-
623 CoV-2 variants. Presented cryo-electron microscopy structure was obtained from
624 Protein Data Bank (accession ID: 6VSB) and processed in UCSF Chimera software.
625 Unique, variant-specific substitutions were indicated in a side view.

626 **Figure 9:** Comparison of amino acid substitution localisation in spike proteins of SARS-
627 CoV-2 variants. Presented cryo-electron microscopy structure was obtained from
628 Protein Data Bank (accession ID: 6VSB) and processed in UCSF Chimera software.
629 Unique, variant-specific substitutions were indicated in a side view.

