1    **Evolution to increased positive charge on the viral spike protein may be part of**

2    **the adaptation of SARS-CoV-2 to human transmission.**

3

4    **Matthew Cotten\*, My V.T. Phan**

5

6    Author affiliations: Medical Research Council–University of Glasgow Centre for Virus

7    Research, Glasgow, Scotland, UK (MC), UK Medical Research Council–Uganda

8    Virus Research Institute and London School of Hygiene and Tropical Medicine

9    Uganda Research Unit, Entebbe, Uganda (MC, MVTP); UK

10

11   *Correspondence to Matthew Cotten, email: matthew.cotten@lshtm.ac.uk

12

13   **Abstract**

14   The severe acute respiratory syndrome coronavirus 2 ( SARS-CoV-2), the causative agent of

15   the coronavirus disease 2019 (COVID-19) pandemic, continues to evolve and infect individuals. The

16   exterior surface of the SARS-CoV-2 virion is dominated by the spike protein and the current work

17   examined spike protein biochemical features that have changed during the 2 years that SARS-CoV-2

18   has infected humans. These biochemical properties may influence virion survival and promote

19   movement through the environment and within the human airway to reach target cells to bind, enter

20   and establish the next round of infection.  In addition to selective pressure to avoid immune

21   recognition of viral proteins, we hypothesised that SARS-CoV-2 emerged from an animal reservoir

22   capable of human infection and transmission but in a sub-optimum state and a second level of

23   selective pressure is acting on these biochemical features. Our analysis identified a striking change in

24   spike protein charge, from -8.3 in the original Lineage A and B viruses to -1.26 in the current Omicron

25   viruses. In summary, we conclude that in addition to immune selection pressure, the evolution of

26   SARS-CoV-2 has also altered viral spike protein biochemical properties. Future vaccine and

27   therapeutic development should also exploit and target these biochemical properties.

28

29   **Introduction**

30   The severe acute respiratory syndrome coronavirus 2 ( SARS-CoV-2), the causative agent of

31   the coronavirus disease 2019 (COVID-19) epidemic, continues to evolve and infect individuals.

32   Similar to other viruses, the SARS-CoV-2 virion biochemical properties play an important role in

33   controlling virus transmission. After replication in an infected individual and release from an infected

34   cell, onward transmission requires survival of the virion to reach susceptible cells in a new host

35   individual initiating the next round of infection. The physical properties of the surface proteins of the

36   virus such as charge, size, hydrophobicity and folding may influence movement of the virion through

37  the environment, promoting or limiting binding of the virion to the external surfaces. Once reaching a

38  susceptible individual, virion physical properties may influence movement within the human airway

39  and determine the ability of an infecting virion to reach target cells to bind, enter and replicate

40  (Adamczyk et al. 2021). The exterior surface of the SARS-CoV-2 virion is dominated by the spike

41  protein and the current work examines simple spike protein features that have changed during the 2

42  years of the SARS-CoV-2 pandemic. In addition to selective pressure to avoid immune recognition of

43  viral proteins, we hypothesise that SARS-CoV-2 emerged from an animal reservoir capable of human

44  infection and transmission but in a sub-optimum state. Additionally,  there is a second level of

45  selective pressure to adjust to the physical transmission between humans. Evidence for this

46  adaptation can be found in changes in the SARS-CoV-2 spike protein over recent evolution. With over

47  11 million SARS-CoV-2 genomic sequences generated globally from across the pandemic, many of

48  these sequences have intact spike gene sequences that can be used to monitor change across the 2

49  years of human host evolution of this virus.

50          Much of the observed spike protein substitutions may be in response to the developing

51  immune response to this new pathogen, which is reflected in substitutions occurring in the immune-

52  exposed S1 domain of the spike protein and there is ample evidence that many of these spike protein

53  changes allow escape from host immunity (Tzou et al. 2022)(Greaney, Loes, et al. 2021)(Greaney,

54  Starr, et al. 2021)(Greaney et al. 2022) (Cao et al. 2022) (Dejnirattisai et al. 2022) (DeGrace et al.

55  2022). There may also be evolutionary selection for protein changes that improve host interactions

56  apart from immune evasion. These include altering spike/receptor binding kinetics, protease cleavage

57  events, tertiary structure (S1/S2 interactions after cleavage) or the physical properties of the virion

58  (charge, hydrophobicity, and protein folding or secondary structure) in ways that might improve

59  transmission. To explore the role of the biochemical features of the spike protein in human

60  transmission, we monitored changes in spike biochemical features over the two years that SARS-

61  CoV-2 has been evolving in humans and report an increase in spike protein positively charge

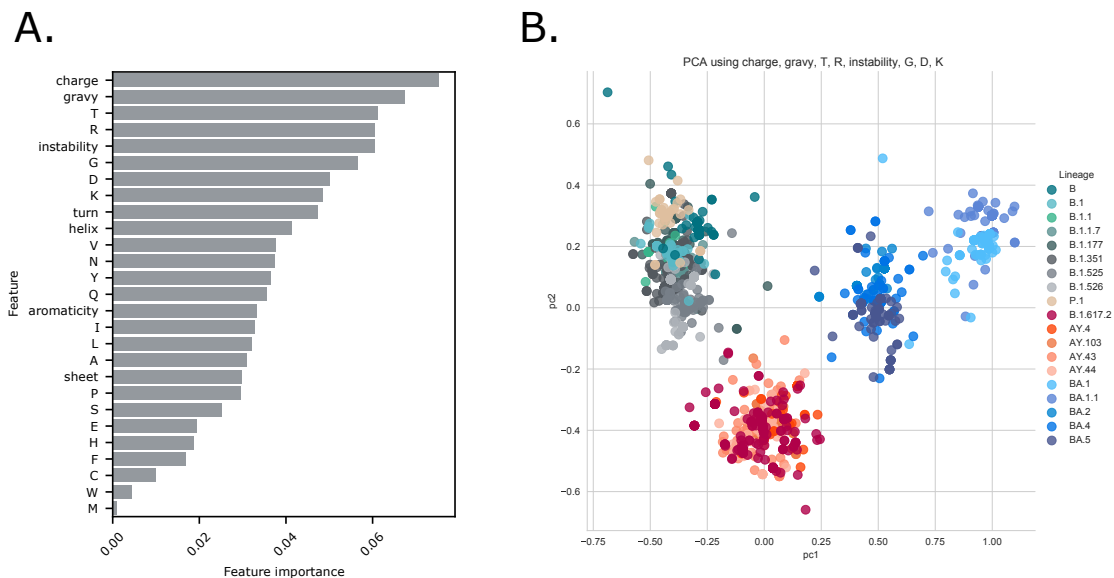62  especially among the virus lineages that were highly prevalent.

63

64  **Results.**

65          The SARS-CoV-2 spike protein physical features were calculated from spike protein

66  sequences from across 2 years of the COVID-19 epidemic. Features that could be quantitated from

67  protein sequence were used (see Methods), including charge at pH 7.4, Kyle and Doolittle GRAVY

68  score (Kyte & Doolittle 1982) (which is a measure of hydrophobicity), an instability index derived from

69  dipeptide content (Guruprasad et al. 1990), properties influencing protein folding (percent helix, fold or

70  sheet as predicted from amino acid content), individual amino acid total fraction and di-amino acid

71  total fraction.

72          A dominant pattern of SARS-CoV-2 evolution during the two years of human adaptation has

73  been the regular appearance and the subsequent regional and then global dominance of lineages.

74  These lineages typically encode a small set of amino acid changes from earlier lineages, many of
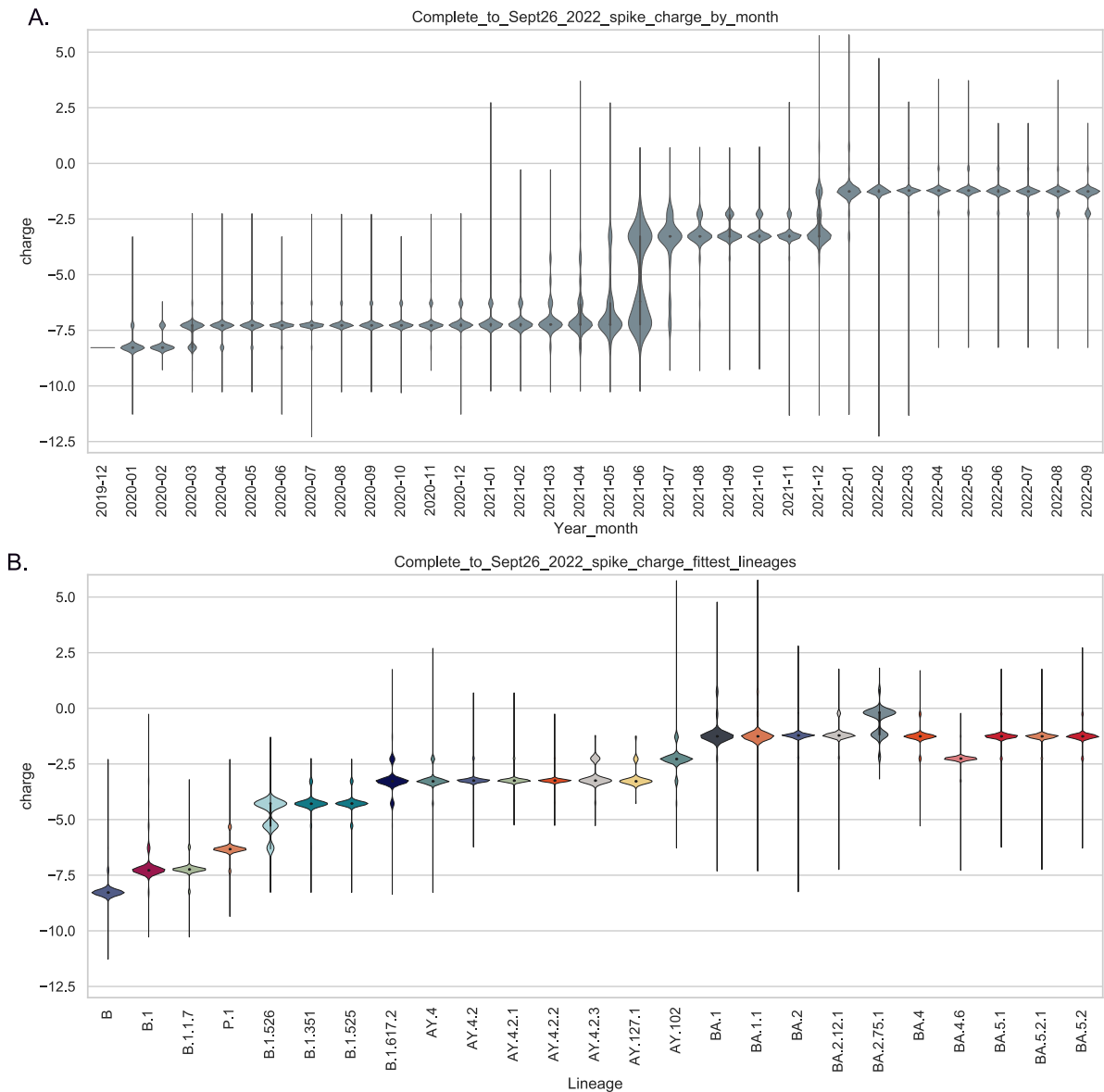
75    which are likely to provide temporary or long-term advantage for the viral lineage. An analysis was

76    performed to identify spike physical features most strongly linked with SARS-CoV-2 lineages (Figure

77    1). The first 300 reported genomes from each major lineage were collected, spike protein sequences

78    were extracted and the physical features of each protein were collected into a matrix. The top features

79    distinguishing SARS-CoV-2 lineages were identified with charge as the most important feature (Figure

80    1A). A principal component analysis using the top 8 features (charge, gravy, fraction T, fraction R,

81    instability, fraction G, fraction D and fraction K), provided clustering of spike sequences by lineage

82    (Figure 1B). These results support the idea that spike protein charge (among other features) is an

83    important determinant of the lineages that have evolved during the first two years of the COVID-19

84    epidemic.



85
86 **Figure 1. Identification of spike protein charge association with SARS-CoV-2 lineage**. **Panel A:**

87    A set of 300 spikes sequences extracted from the first 300 SARS-CoV-2 genomes per lineage (by

88    date of collection) was analyzed, features for each sequence were collected (see Methods). SKLearn

89    feature selection (Pedregosa, F. and Varoquaux, G. and Gramfort, A. and Michel, V. et al. 2011) was

90    used to identify features that most accurately identified the sequence lineage. The importance of

91    features were ranked in order. **Panel B:** The top 8 features (charge, gravy, fraction T, fraction R,

92    instability, fraction G, fraction D, fraction K) were further used in a principal component analysis to

93    cluster the same set of SARS-CoV-2 spike sequences. Each node represents a single spike

94    sequence, nodes were coloured by Pangolin lineage assigned to the genome from which the spike

95    sequence was obtained. Lineage colouring is explained in the figure legend to the right**.**

96
97
98       Changes in charge of spike protein across the epidemic were investigated. Plotting total spike

99    charge for all genomes per month of the epidemic showed a clear pattern of increase in charge over

100    two years of evolution (Figure 2, panel A). Median spike charge was -8.3 in the original SARS-CoV-2

101     viruses reported in late 2019 to early 2020, by March 2020, an increase in positive charge to -7.28

102     was observed. Subsequently, an additional increase in positive charge occurred in mid-2021 to -3.28,

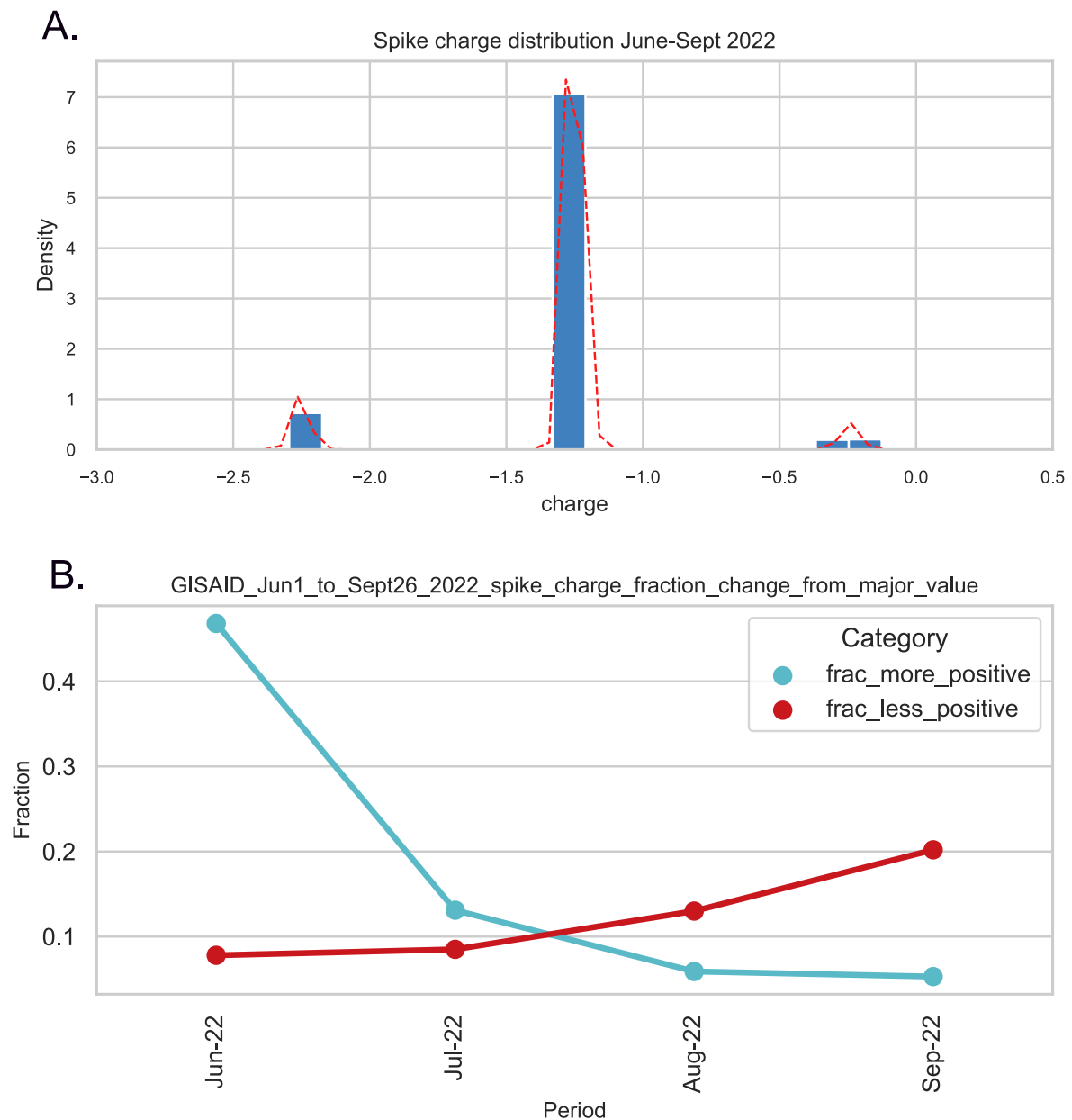103     and most recently a charge increase occurred in late 2020/early 2021 to -1.26.

104



105
106     **Figure 2**.**Panel A. Total SARS-CoV-2 spike charge per epidemic month**. All available SARS-CoV-

107     2 genomes up to 24 September 2022 were retrieved from GISAID (GISAID 2020) and the spike

108     protein sequence was extracted (if intact). Total charge at pH 7.4 was calculated and values were

109     plotted using a violin plot by month of sample collection. For each epidemic month the violin plot

110     depicts the distributions of calculated spike charge for all available SARS-CoV-2 genomes. **Panel B.**

111     **Spike charge in major SARS-CoV-2 lineages.** For each lineage, all available spike sequences were

112     collected (up to 24 September 2022), total charge was measured and violin plots prepared to show

113     the charge distribution by lineage. Lineages (indicated at bottom of chart) were ordered by their

114     appearance in the epidemic.

4

115

116    These spike protein charge increases can be attributed to the major successful lineages

117    reported over time (Figure 2B). The B.1, B.1.1 and B.1.1.7 (Alpha) lineages that dominated the first

118    year of the epidemic encoded spike proteins with charges between -8 and -6 while the B.1.351 (Beta)

119    and B.1.525 (Eta) lineages showed a further increase in charge to around -4.5. The B.1.617.2 (Delta)

120    lineage and sublineages (AY.x) displayed further increase in charge. Most recently, the Omicron

121    variants (including BA.1, BA.1.1, BA.2, BA.2.12,BA.3, BA.4 and BA.5) show further spike charge

122    increases with the majority of Omicron encoded spike proteins showing charge at -1.26 (Figure 2,

123    panel B).

124    Some indication of functional consequences of the observed changes in spike charge can be

125    obtained from the location on the charged amino acid substitutions in the spike protein. Sets of spike

126    sequences (extracted from the first 300 reported genomes per select lineage) were processed to

127    illustrate the changes to more negative charge (blue) or more positive charge (orange/red) in the

128    protein relative to the initial Lineage B genome sequences (Supplementary Figure 1). The initial

129    change in charge was a substitution of  an aspartic acid residue (D, with a calculated charge of -1) by

130    a glycine (G, neutral). In some early lineages (e.g. A.23.1), proline (P) at position 681was substituted

131    with the positively charged arginine (R),or Q680 was substituted with a partially charged histidine H

132    residue. The P681R positive substitution promotes furin cleavage and activation of the spike protein

133    for cell fusion (Lubinski et al. 2022)(Liu et al. 2022). The Delta lineage spike proteins encoded

134    additional positive charge in the ACE2 binding region, as well as in the far amino terminal region and

135    near the heptad repeat (HR1) which may also enhance membrane fusion activity. More recently, a

136    number of positive substitutions have occurred in the Omicron lineage virus spike proteins with

137    predominance of positively charged changes in the receptor binding domain (Supplementary Figure

138    1), suggesting a role of increased charge in spike/receptor interactions.

139    It is probable that the spike protein has an upper limit to the amino acid charge that it can

140    allow for proper folding, assembly and function. This upper charge value will be determined by the

141    acquisition of optimum transmission properties in balance with immune selection. After the regular

142    increase of spike protein charge observed up to the appearance of the Omicron lineages,  an

143    indication of a stasis in positively-charged amino acid accumulation is now displayed by SARS-CoV-2

144    Omicron lineages. The majority of Omicron sub-lineages remain at spike charge -1.26 (Figure 3A)

145    although a few specific Omicron sub-lineages show changes toward more positive or negative charge

146    (e.g. BA.2.75.1 more positive, BA.4.6  more negative, as illustrated in Figure 2B) with the additional

147    changes often associated with immune selection. To monitor the current trends of spike protein

148    changes, we calculated the fraction of reported genomes with spike charge greater than or less than

149    the Omicron mean charge of -1.26 and documented how these fractions had changed over the last 4

150    months of the pandemic (Figure 3B). The majority of encoded spike proteins are almost exclusively

151    from Omicron lineage viruses and show a charge of -1.26. However, a small fraction of genomes

152    encode spikes proteins with slightly more or less charge (Figure 3B) with the greater trend (20% of all

153    reported genomes in September 2022)  showing more negative charge (Figure 3B).

154

**A.**



**B.**



155

156    Figure 3. Recent changes in spike protein charge. **Panel A**: All available spike proteins from genomes

157    with sample collection dates of June-Sept 2022 were analyzed for total spike charge. A histogram of

158    the calculate total spike charges for the entire set is shown in panel A with the **kernel density**

159    **estimation** (**KDE**) line in red.  A major peak at -1.26 is observed wit small outlier peaks of genomes

160    with more negative and more positive spike proteins **Panel B**: For each month (over the period June 1

161    to Sept 26 2022) the fraction of reported genomes for that month with charge greater than or less than

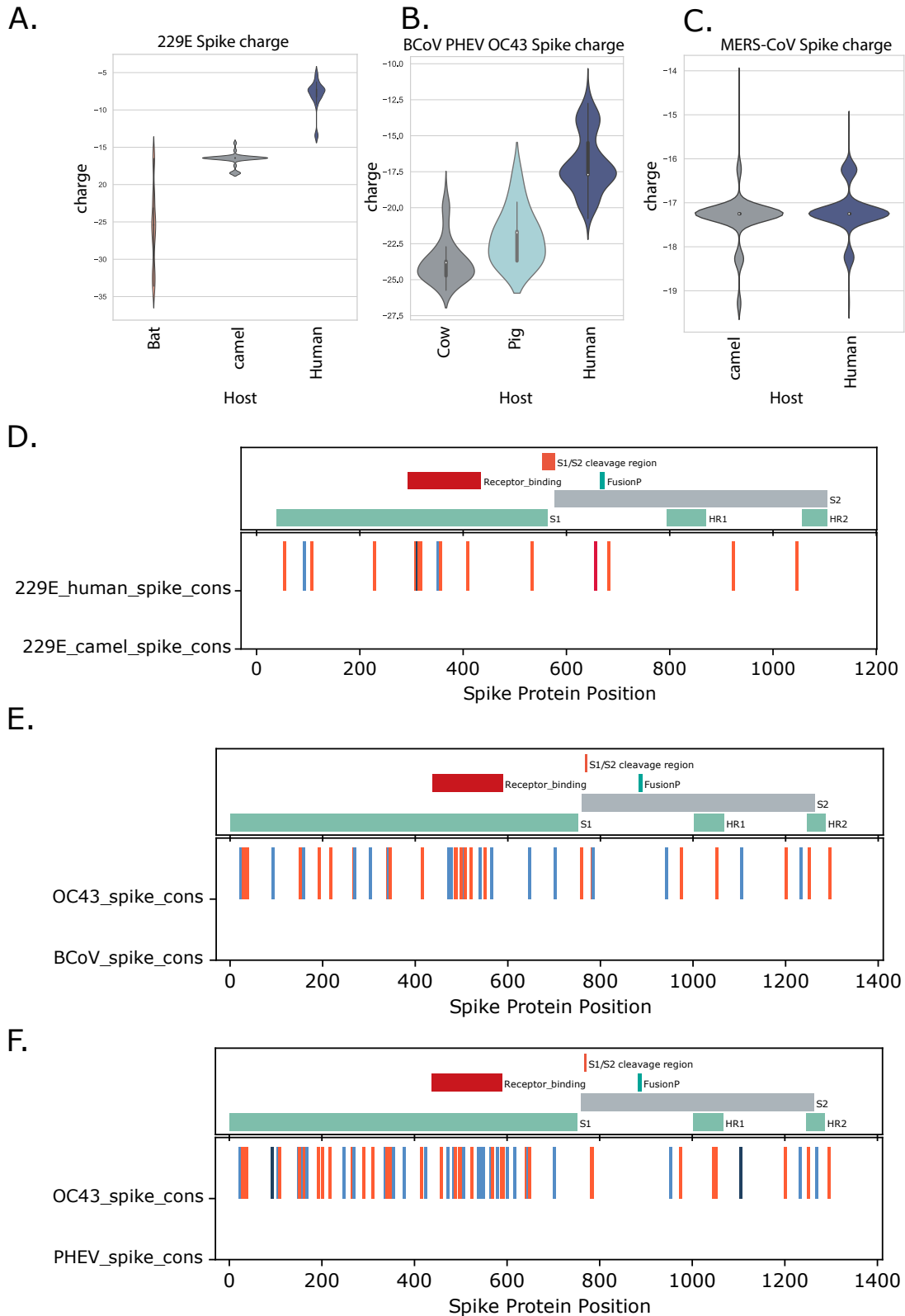162    the majority value of -1.26 was calculated.

6

163    Lastly, we investigated if a similar pattern of spike charge evolution could be observed in

164    other coronaviruses that have made a transition to human transmission. In recent history, several

165    coronaviruses (in addition to SARS-CoV-2) have been observed to jump hosts. For example,

166    coronavirus 229E is commonly detected in humans and very close coronaviruses have been identified

167    in bats (Victor Max Corman et al. 2015) (Tao et al. 2017) and camels (Victor M. Corman et al. 2016)

168    (Sabir et al. 2016) suggesting movement of the virus between hosts.  All available coronavirus 229E

169    full genomes sequences were retrieved from GenBank, the spike coding region was extracted from

170    the genomes, translated and total charge was calculated. A difference from -26 to -8, or almost 18

171    charge units is seen comparing 229E-like viruses from bats to 229E from humans (Figure 4A)  and

172    almost 9 charge unit difference  was observed in spike median charge comparing 229E viruses from

173    camel vs. human infections (Figure 4A).

174    Infection with coronavirus OC43 is common in humans and closely related viruses are found

175    in cattle (bovine coronavirus, BCoV)  (Vijgen et al. 2006)  and pigs (porcine hemagglutinating

176    encephalomyelitis virus, PHEV) (Vijgen et al. 2005) (Vijgen et al. 2006). Comparing the three OC43-

177    type virus groups, the human virus OC43 has an increased charge of ca. 5 units compared to PHEV

178    and ca.7 units compared to BCoV (Figure 4B).

179    The commonly known host for the Middle East Respiratory Syndrome coronavirus (MERS-

180    CoV) is dromedary camels; however, zoonosis and serious human infections occur frequently

181     (Cotten et al. 2013) (Cotten et al. 2014) (Memish et al. 2014) (Zhou et al. 2021) (So et al. 2019) as

182    reviewed in (Peiris & Perlman 2022).  From 698 full MERS-CoV genomes available in GenBank, there

183    was no strong difference in the encoded spike charge of virus sequences derived from human versus

184    camels infections (Figure 4C).

185    Considering the location of the charge differences in the spike proteins, for coronavirus 229E,

186    the charge increases occurred throughout the protein, although there is a slightly higher number of

187    positive changes in the receptor binding region of the human infection derived viruses (Figure 4D).

188    For OC43, the porcine and human viruses also show increases in positive charge throughout the

189    spike protein, the porcine PHEV also showed a slight enrichment in positive charge in the receptor

190    binding region (Figure 4E).

191

**Figure 4. Spike charges from select groups of coronaviruses that have moved into humans**

**(Panels A-C)**. All available full genomes for the indicated coronaviruses were retrieved from

196   GenBank, the spike coding region was identified and translated into protein and total charge at ph 7.4

197   was calculated. Violin plots indicate the charges of each collection of spike proteins, median values

198   are indicated by the open square. **Panel A**. Coronavirus 229E from bat, camel or human infections,

199   **Panel B**. BCoV (from bovine infections) PHEV (from porcine infections) and OC43 (from human

200   infection), **Panel C.** MERS-CoV from camel or human infection. **Panel D-F: Consensus spike**

201   **protein sequences were generated from the indicated virus groups and charged amino acid**

202   **changes were determined.** Charge changes were colored from dark blue (change from positively to

203   negatively charged amino acid (AA)), blue change from neutral to negatively charged AA), orange

204   (change from neutral to positively charged AA) and red (change from negative to positively charged

205   AA). **Panel D:** 229E spike from human infections compared to 229E spike from camel infections,

206   **Panel E**: Human OC43 spike compared to BCoV spike, **Panel F**: Human OC43 spike compared to

207   PHEV spike. Key spike protein features of each group's spike protein are shown in the upper portion

208   of each panel.

209

210   **Discussion**

211         After more than two years of the COVID-19 pandemic and with the availability of >11 million

212   SARS-CoV-2 genome sequences, a trend of SARS-CoV-2 spike protein charge can be observed,

213   with successive lineages showing an increase in positive charge over earlier lineages. Over the

214   course of the pandemic, the SARS-CoV-2 spike protein has evolved from a protein with a total charge

215   of -8.28 in the original Lineage A and B viruses to a protein with a total charge of -1.26 in the majority

216   of the currently circulating Omicon lineage viruses. This pattern has been noted previously

217   (Pawłowski 2021) (Nie et al. 2022). We expand on these observations, and document lineage

218   patterns and sites of change in the spike protein and explore similar phenomena of evolution to more

219   positive charge in two other coronaviruses (OC43 and 229E) that have moved between animals and

220   humans.

221         This study does not identify a mechanistic basis for the increased spike charge although there

222   are several possible transmission steps that might be promoted by increasing charge. Exposed,

223   positively charged spike amino acids should promote interactions with negatively charged cellular

224   structures. Interactions with negatively-charged heparin have been reported with SARS-CoV-2 spike

225   (Kim et al. 2020) and negatively-charged sialylated glycans are reported to promote entry of SARS-

226   CoV-2 (Nguyen et al. 2022). The upper respiratory tract is coated by and protected by mucins,

227   frequently modified with sialic acid or phosphorylated, high mannose N-glycans (Byrd-Leotis et al.

228   2021) which present a negatively charged matrix that could either promote or protect against viral

229   transmission. The SARS-CoV-2, OC43, and BCoV virions display binding to negatively charged

230   carbohydrate structures found in the airway (Byrd-Leotis et al. 2021) and the ionic environment of the

231   human upper respiratory tract may favour binding and transmission of viruses with increased positive

232   charge. Perhaps it is not surprising that both OC43 and 229E coronaviruses exhibited increases in

233   spike positive charge after moving from animal hosts (cow, pig and camel) to human hosts (Figure 5).

234   A similar change in MERS-CoV was not observed, however MERS-CoV currently shows only limited

235   human to human transmission with most known transmission chains ending after 2 to 3 human to

236   human transmission events as shown in (Assiri et al. 2013) and (Cotten et al. 2013). MERS-CoV

237   might not have experienced sufficient number of human replication cycles or have undergone the

238   same level of selection for human transmission that OC43, 229E and SARS-CoV-2 have experienced.

239   For both OC43 and 229E coronaviruses moving to humans, the broad location of the positive

240   changes across the spike protein sequence suggested that positive charge may be promoting several

241   functions including receptor binding, furin cleavage, cell fusion as well as antigenic changes or less

242   specific changes to avoid or promote ionic interactions during transmission.

243       There is likely a limit to the accumulation of positively charged residues in the SARS-CoV-2

244   spike protein. Functional constraints exist, there may also be penalties associated with non-specific

245   binding due to excess positive charge, and there are certainly charge influences on protein folding

246   and higher order protein interactions (Creighton 2002). Our prediction is that the SARS-CoV-2 protein

247   will reach some upper limit of charge defined by these constraints. Indeed, we observe that the

248   majority of Omicron lineages encode spike proteins with charge -1.26, after more than 6 months of

249   evolution (Figures 2A and 2B). A small fraction of genomes with more positive charge or less positive

250   charge have appeared, but the global tendency across all reported genomes from June to September

251   2022 is a modest decline in the positive charge (Figure 3b) which suggest the upper limit to charge

252   has been reached.

253       Could these changes in spike charge have occurred by chance and not be a response to

254   selective pressure? Of the 20 standard amino acids (AA), only 2 AA have negatively charged side

255   chains, 2 AA have positively charged side chains while the remaining 16 AA are neutral at pH 7.4..

256   Assuming equal probability of any AA change, there is an 18/20 chance of a negative AA being

257   substituted by a neutral or positively charged AA and the majority of change opportunities would result

258   in loss of negative charge. However, natural selection is more complex, because the genetic code

259   uses 3 adjacent nucleotides to encode an AA, there are multiple encoding possibilities for each AA,

260   the codon redundancy is not identical for each AA and the number of nucleotide changes required to

261   produce any particular AA change can be 1, 2 or 3. This has resulted in an evolved protein stability in

262   the genetic code (Chan et al. 2020) with AA changes that maintain rather than change physical

263   properties (negative, positive, polar, non-polar, aromatic) more likely based on the codon array

264   (Livingstone & Barton 1993) and the nucleotide changes required for an AA change. For example, the

265   probability of a negative AA to negative AA change is 0.333 while the probabilities of change of a

266   negative AA to a non-polar, aromatic, polar or positive AA are 0.051, 0.044, 0.028 and 0.044

267   respectively, with changes away from a negative charged AA nearly 10-fold less likely to occur than

268   conserving the negative charge at that position (Livingstone & Barton 1993).  For these reasons, it

269   appears that the accumulation of positive charge on spike protein has not occurred by chance and is

270   likely providing some selective advantage for the virus. It should also be noted that the observed

271   charge changes in exposed virion proteins seem to be limited to spike. Two additional SARS-CoV-2

10

272    proteins are externally exposed, the E protein (ORF4) and the M protein (ORF5), showing no

273    consistent change in the charge of either of these proteins across the 2 years of the epidemic (results

274    not shown).

275         Obermeyer et al. documented AA substitutions associated with SARS-CoV-2 fitness

276    (Obermeyer et al. 2022). Consistent with the idea that the increase in positive charges is not by

277    chance, of the top 20 substitutions increasing SARS-CoV-2 fitness, 14 substitutions were in the spike

278    protein, among which 4 were changes that increased positive charge while only 1 of 14 introduced a

279    negative charge in spike (Obermeyer et al. 2022).

280         Natural selection could be acting on multiple features of the spike protein. The necessity to

281    avoid host immune responses is likely to be the major selective force acting on the virus. This results

282    in the amino acid changes, which in turn are determined by epitopes. The selection for increased

283    charge in the spike protein is probably occurring in the background, not as a major shift needed to

284    bypass immune responses. However, the increase in charge may improving survival and transmission

285    in humans in subtle ways, and this advantage, when multiplied over the millions of infections can

286    provide some of the growth and infection advantages seen by new SARS-CoV-2 variants. It is

287    proposed that the N764K, N856K and N969K substitutions (all increasing spike positive charge) may

288    enhance S1/S2 subunit interactions after proteolytic processing of the spike protein, resulting in

289    reduced S1 shedding and improving transmission (Martin et al. 2022) Increased charge may also alter

290    receptor interactions. In the Omicron (BA.1) spike protein, the Q493R and Q498R substitutions are

291    predicted to allow two additional salt bridges with ACE2 receptor position 35Glu and 38Glu (McCallum

292    et al. 2022). Indeed, looking at the timing of charge shifts in each major lineage, the changes to more

293    positive charge accumulate later than the changes that first allow a lineage to emerge and dominate

294    global infections. In this model, the primary spike changes are driven by immune selection and allow a

295    new lineage to bypass existing immune responses. Once a successful new variant emerges, the large

296    number of new infections allow selection for the accumulation of beneficial positive charge changes.

297    The similar pattern of increased positivity of spike protein in other coronaviruses that have moved

298    between animals and humans (OC43, 229E, Figure 4) suggest that the change in surface protein

299    charge may be a more general phenomenon with coronaviruses and might be a useful parameter to

300    examine when monitoring zoonosis.  This study provides a framework to monitor viral evolution

301    through changes in biochemical properties, which can be easily applied to other viruses important to

302    public and global health. An important note, our analyses on viral spike protein biochemical properties

303    to monitor virus evolution are not meant to replace traditional phylogenetic analyses. The observed

304    pattern of biochemical properties changes should completement phylogenetic signals. However, in

305    situations where there are limited sequences available to produce reliable phylogenetic signals (e.g.

306    the 229E and OC43 viruses examined in Figure 4), this kind of analysis using virus biochemical

307    properties from different host species would certainly help provide important information on the virus

308    evolution, zoonosis as well as aiding the prediction of patterns of viral changes.

309    In conclusion, our study provides an novel analytical framework to monitor viral evolution

310    through changes in biochemical properties, which can be easily applied to other viruses important to

311    public and global health. We also showed that natural virus evolution is more complicated and may

312    involve multiple factors including immune selection, as well as spike protein biochemical properties.

313    The observation of increase of SARS-CoV-2 spike protein charge over time provides useful

314    information for future vaccine and therapeutic development.

315

316    **Methods**

317    Full alignments of SARS-CoV-2 genomes were obtained from GISAID (GISAID 2020) with

318    collection dates to 15 June 2022. All spaces in fasta IDs were removed using sed (sed -i -e 's/ /_/g'

319    msa_xxxx.fasta), the alignment was dealigned ("-" characters removed) and genomes were classified

320    using Pangolin (Áine O'Toole et al. 2020) with the most recent database updates (pangolin v4.1.1,

321    pangolin-data v1.11

322    constellations v0.1.10 and scorpio v0.3.17). The spike coding region from each genome (if present

323    and intact (no Ns)) was translated into protein. Features of the protein that could be quantitated from

324    the spike protein sequence were determined using the ProteinAnalysis functions from BioPython

325    (Cock et al. 2009). These features included charge at pH 7.4, Kyle and Doolittle GRAVY score

326    (Kyte & Doolittle 1982) (a measure of hydrophobicity), an instability index derived from dipeptide

327    content (Guruprasad et al. 1990), the total percent helix, fold or sheet properties of the protein and the

328    total fractions of individual amino acids and fractions of di-amino acids. A matrix of all spike protein

329    features plus collection date, and lineage was prepared and used for analysis. Similar analyses were

330    performed for other coronaviruses such as 229E, OC43 and MERS-CoV by retrieving all complete

331    genomes available from GenBank (15 June 2022). The spike protein was also extracted using the

332    same method as aforementioned. Additional details are provided in the figure legends. The python

333    code used for the analyses is available here: https://github.com/mlcotten13/SARS-CoV-
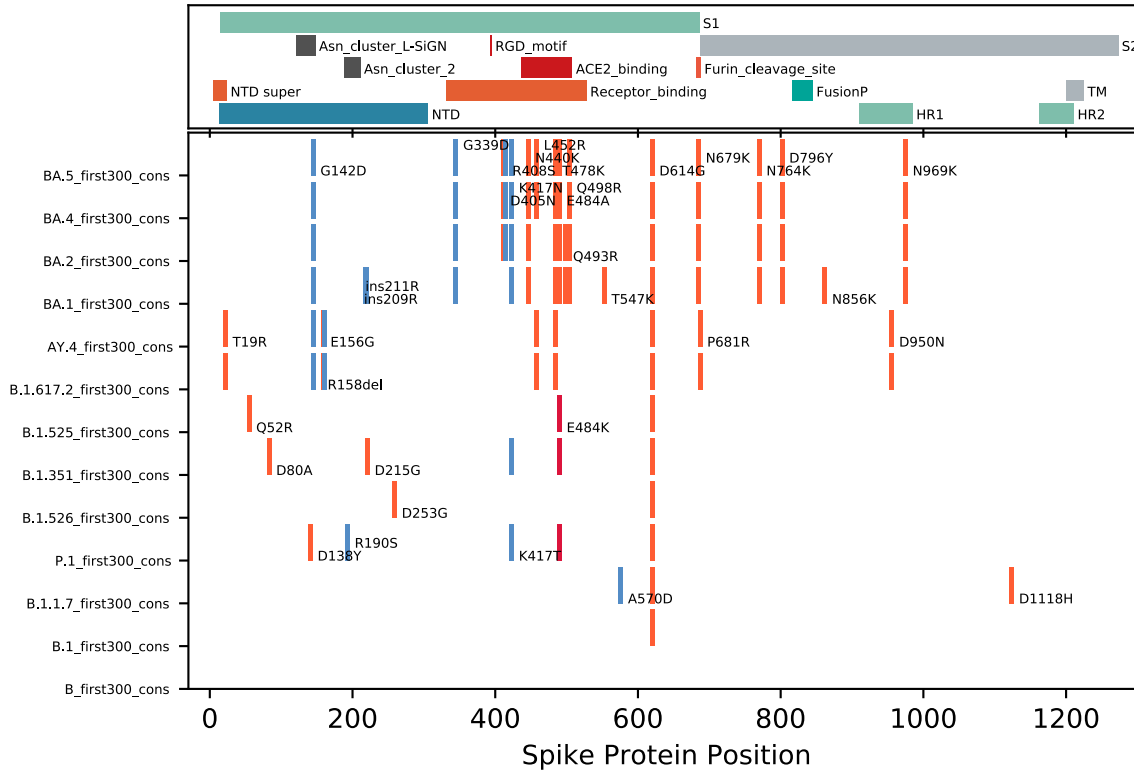
334    2_spike_charge .

335

346

## References

348 Adamczyk, Z., Batys, P., & Barbasz, J. (2021). 'SARS-CoV-2 virion physicochemical characteristics
349   pertinent to abiotic substrate attachment', *Current Opinion in Colloid & Interface Science*, 55:
350   101466. DOI: 10.1016/j.cocis.2021.101466
351 Áine O'Toole et al. (2020). 'Phylogenetic Assignment of Named Global Outbreak LINeages
352   (PANGOLIN)',.
353 Assiri, A., McGeer, A., Perl, T. M., Price, C. S., Al Rabeeah, A. A., Cummings, D. A. T., Alabdullatif, Z.
354   N., et al. (2013). 'Hospital Outbreak of Middle East Respiratory Syndrome Coronavirus', *New*
355   *England Journal of Medicine*, 369/5: 407–16. DOI: 10.1056/NEJMoa1306742
356 Byrd-Leotis, L., Lasanajak, Y., Bowen, T., Baker, K., Song, X., Suthar, M. S., Cummings, R. D., et al.
357   (2021). 'SARS-CoV-2 and other coronaviruses bind to phosphorylated glycans from the
358   human lung', *Virology*, 562: 142–8. DOI: 10.1016/j.virol.2021.07.012
359 Cao, Y., Wang, J., Jian, F., Xiao, T., Song, W., Yisimayi, A., Huang, W., et al. (2022). 'Omicron
360   escapes the majority of existing SARS-CoV-2 neutralizing antibodies', *Nature*, 602/7898:
361   657–63. DOI: 10.1038/s41586-021-04385-3
362 Chan, K.-F., Koukouravas, S., Yeo, J. Y., Koh, D. W.-S., & Gan, S. K.-E. (2020). 'Probability of
363   change in life: Amino acid changes in single nucleotide substitutions', *Biosystems*, 193–194:
364   104135. DOI: 10.1016/j.biosystems.2020.104135
365 Cock, P. J. A., Antao, T., Chang, J. T., Chapman, B. A., Cox, C. J., Dalke, A., Friedberg, I., et al.
366   (2009). 'Biopython: freely available Python tools for computational molecular biology and
367   bioinformatics', *Bioinformatics*, 25/11: 1422–3. DOI: 10.1093/bioinformatics/btp163
368 Corman, Victor M., Eckerle, I., Memish, Z. A., Liljander, A. M., Dijkman, R., Jonsdottir, H., Juma
369   Ngeiywa, K. J. Z., et al. (2016). 'Link of a ubiquitous human coronavirus to dromedary
370   camels', *Proceedings of the National Academy of Sciences*, 113/35: 9864–9. DOI:
371   10.1073/pnas.1604472113
372 Corman, Victor Max, Baldwin, H. J., Tateno, A. F., Zerbinati, R. M., Annan, A., Owusu, M., Nkrumah,
373   E. E., et al. (2015). 'Evidence for an Ancestral Association of Human Coronavirus 229E with
374   Bats', (S. Schultz-Cherry, Ed.)*Journal of Virology*, 89/23: 11858–70. DOI: 10.1128/JVI.01755-
375   15
376 Cotten, M., Watson, S. J., Kellam, P., Al-Rabeeah, A. A., Makhdoom, H. Q., Assiri, A., Al-Tawfiq, J.
377   A., et al. (2013). 'Transmission and evolution of the Middle East respiratory syndrome
378   coronavirus in Saudi Arabia: a descriptive genomic study', *Lancet (London, England)*,
379   382/9909: 1993–2002. DOI: 10.1016/S0140-6736(13)61887-5
380 Cotten, M., Watson, S. J., Zumla, A. I., Makhdoom, H. Q., Palser, A. L., Ong, S. H., Al Rabeeah, A.
381   A., et al. (2014). 'Spread, circulation, and evolution of the Middle East respiratory syndrome
382   coronavirus', *mBio*, 5/1: e01062-13. DOI: 10.1128/mBio.01062-13
383 Creighton, T. E. (2002). *Proteins: structures and molecular properties*., 2. ed., [Nachdr.]. New York:
384   Freeman.
385 DeGrace, M. M., Ghedin, E., Frieman, M. B., Krammer, F., Grifoni, A., Alisoltani, A., Alter, G., et al.
386   (2022). 'Defining the risk of SARS-CoV-2 variants on immune protection', *Nature*, 605/7911:
387   640–52. DOI: 10.1038/s41586-022-04690-5
388 Dejnirattisai, W., Huo, J., Zhou, D., Zahradník, J., Supasa, P., Liu, C., Duyvesteyn, H. M. E., et al.
389   (2022). 'SARS-CoV-2 Omicron-B.1.1.529 leads to widespread escape from neutralizing
390   antibody responses', *Cell*, 185/3: 467-484.e15. DOI: 10.1016/j.cell.2021.12.046
391 GISAID. (2020). 'The GISAID Initiative',.
392 Greaney, A. J., Loes, A. N., Crawford, K. H. D., Starr, T. N., Malone, K. D., Chu, H. Y., & Bloom, J. D.
393   (2021). 'Comprehensive mapping of mutations in the SARS-CoV-2 receptor-binding domain
394   that affect recognition by polyclonal human plasma antibodies', *Cell Host & Microbe*, 29/3:
395   463-476.e6. DOI: 10.1016/j.chom.2021.02.003
396 Greaney, A. J., Starr, T. N., Barnes, C. O., Weisblum, Y., Schmidt, F., Caskey, M., Gaebler, C., et al.
397   (2021). 'Mapping mutations to the SARS-CoV-2 RBD that escape binding by different classes
398   of antibodies', *Nature Communications*, 12/1: 4196. DOI: 10.1038/s41467-021-24435-8

399    Greaney, A. J., Starr, T. N., & Bloom, J. D. (2022). 'An antibody-escape estimator for mutations to the
400           SARS-CoV-2 receptor-binding domain', *Virus Evolution*, 8/1: veac021. DOI:
401           10.1093/ve/veac021
402    Guruprasad, K., Reddy, B. V. B., & Pandit, M. W. (1990). 'Correlation between stability of a protein
403           and its dipeptide composition: a novel approach for predicting *in vivo* stability of a protein from
404           its primary sequence', *Protein Engineering, Design and Selection*, 4/2: 155–61. DOI:
405           10.1093/protein/4.2.155
406    Kim, S. Y., Jin, W., Sood, A., Montgomery, D. W., Grant, O. C., Fuster, M. M., Fu, L., et al. (2020).
407           'Characterization of heparin and severe acute respiratory syndrome-related coronavirus 2
408           (SARS-CoV-2) spike glycoprotein binding interactions', *Antiviral Research*, 181: 104873. DOI:
409           10.1016/j.antiviral.2020.104873
410    Kyte, J., & Doolittle, R. F. (1982). 'A simple method for displaying the hydropathic character of a
411           protein', *Journal of Molecular Biology*, 157/1: 105–32. DOI: 10.1016/0022-2836(82)90515-0
412    Liu, Y., Liu, J., Johnson, B. A., Xia, H., Ku, Z., Schindewolf, C., Widen, S. G., et al. (2022). 'Delta
413           spike P681R mutation enhances SARS-CoV-2 fitness over Alpha variant', *Cell Reports*, 39/7:
414           110829. DOI: 10.1016/j.celrep.2022.110829
415    Livingstone, C. D., & Barton, G. J. (1993). 'Protein sequence alignments: a strategy for the
416           hierarchical analysis of residue conservation', *Bioinformatics*, 9/6: 745–56. DOI:
417           10.1093/bioinformatics/9.6.745
418    Lubinski, B., Frazier, L. E., Phan, M. V. T., Bugembe, D. L., Cunningham, J. L., Tang, T., Daniel, S.,
419           et al. (2022). 'Spike Protein Cleavage-Activation in the Context of the SARS-CoV-2 P681R
420           Mutation: an Analysis from Its First Appearance in Lineage A.23.1 Identified in Uganda',
421           *Microbiology Spectrum*, e0151422. DOI: 10.1128/spectrum.01514-22
422    Martin, D. P., Lytras, S., Lucaci, A. G., Maier, W., Grüning, B., Shank, S. D., Weaver, S., et al. (2022).
423           'Selection analysis identifies clusters of unusual mutational changes in Omicron lineage BA.1
424           that likely impact Spike function', (K. Crandall, Ed.)*Molecular Biology and Evolution*, msac061.
425           DOI: 10.1093/molbev/msac061
426    McCallum, M., Czudnochowski, N., Rosen, L. E., Zepeda, S. K., Bowen, J. E., Walls, A. C., Hauser,
427           K., et al. (2022). 'Structural basis of SARS-CoV-2 Omicron immune evasion and receptor
428           engagement', *Science*, 375/6583: 864–8. DOI: 10.1126/science.abn8652
429    Memish, Z. A., Cotten, M., Meyer, B., Watson, S. J., Alsahafi, A. J., Al Rabeeah, A. A., Corman, V. M.,
430           et al. (2014). 'Human infection with MERS coronavirus after exposure to infected camels,
431           Saudi Arabia, 2013', *Emerging Infectious Diseases*, 20/6: 1012–5. DOI:
432           10.3201/eid2006.140402
433    Nguyen, L., McCord, K. A., Bui, D. T., Bouwman, K. M., Kitova, E. N., Elaish, M., Kumawat, D., et al.
434           (2022). 'Sialic acid-containing glycolipids mediate binding and viral entry of SARS-CoV-2',
435           *Nature Chemical Biology*, 18/1: 81–90. DOI: 10.1038/s41589-021-00924-1
436    Nie, C., Sahoo, A. K., Netz, R. R., Herrmann, A., Ballauff, M., & Haag, R. (2022). 'Charge Matters:
437           Mutations in Omicron Variant Favor Binding to Cells', *ChemBioChem*. DOI:
438           10.1002/cbic.202100681
439    Obermeyer, F., Jankowiak, M., Barkas, N., Schaffner, S. F., Pyle, J. D., Yurkovetskiy, L., Bosso, M.,
440           et al. (2022). 'Analysis of 6.4 million SARS-CoV-2 genomes identifies mutations associated
441           with fitness', *medRxiv: The Preprint Server for Health Sciences*, 2021.09.07.21263228. DOI:
442           10.1101/2021.09.07.21263228
443    Pawłowski, P. (2021). 'SARS-CoV-2 variant Omicron (B.1.1.529) is in a rising trend of mutations
444           increasing the positive electric charge in crucial regions of spike protein S.', *Acta Biochimica
445           Polonica*. DOI: 10.18388/abp.2020_6072
446    Pedregosa, F. and Varoquaux, G. and Gramfort, A. and Michel, V., and Thirion, B. and Grisel, O. and
447           Blondel, M. and Prettenhofer, P., and Weiss, R. and Dubourg, V. and Vanderplas, J. and
448           Passos, A. and, & Cournapeau, D. and Brucher, M. and Perrot, M. and Duchesnay, E. (2011).
449           'Scikit-learn: Machine Learning in Python', *Journal of Machine Learning Research*, 12: 2825–
450           30.
451    Peiris, M., & Perlman, S. (2022). 'Unresolved questions in the zoonotic transmission of MERS',
452           *Current Opinion in Virology*, 52: 258–64. DOI: 10.1016/j.coviro.2021.12.013

453 Sabir, J. S. M., Lam, T. T.-Y., Ahmed, M. M. M., Li, L., Shen, Y., E. M. Abo-Aba, S., Qureshi, M. I., et
454   al. (2016). 'Co-circulation of three camel coronavirus species and recombination of MERS-
455   CoVs in Saudi Arabia', *Science*, 351/6268: 81–4. DOI: 10.1126/science.aac8608
456 So, R. T. Y., Chu, D. K. W., Miguel, E., Perera, R. A. P. M., Oladipo, J. O., Fassi-Fihri, O., Aylet, G., et
457   al. (2019). 'Diversity of Dromedary Camel Coronavirus HKU23 in African Camels Revealed
458   Multiple Recombination Events among Closely Related Betacoronaviruses of the Subgenus
459   Embecovirus', *Journal of Virology*, 93/23: e01236-19. DOI: 10.1128/JVI.01236-19
460 Tao, Y., Shi, M., Chommanard, C., Queen, K., Zhang, J., Markotter, W., Kuzmin, I. V., et al. (2017).
461   'Surveillance of Bat Coronaviruses in Kenya Identifies Relatives of Human Coronaviruses
462   NL63 and 229E and Their Recombination History', (S. Perlman, Ed.)*Journal of Virology*, 91/5:
463   e01953-16. DOI: 10.1128/JVI.01953-16
464 Tzou, P. L., Tao, K., Pond, S. L. K., & Shafer, R. W. (2022). 'Coronavirus Resistance Database (CoV-
465   RDB): SARS-CoV-2 susceptibility to monoclonal antibodies, convalescent plasma, and
466   plasma from vaccinated persons', (J. Bhattacharya, Ed.)*PLOS ONE*, 17/3: e0261045. DOI:
467   10.1371/journal.pone.0261045
468 Vijgen, L., Keyaerts, E., Lemey, P., Maes, P., Van Reeth, K., Nauwynck, H., Pensaert, M., et al.
469   (2006). 'Evolutionary history of the closely related group 2 coronaviruses: porcine
470   hemagglutinating encephalomyelitis virus, bovine coronavirus, and human coronavirus
471   OC43', *Journal of Virology*, 80/14: 7270–4. DOI: 10.1128/JVI.02675-05
472 Vijgen, L., Keyaerts, E., Moës, E., Thoelen, I., Wollants, E., Lemey, P., Vandamme, A.-M., et al.
473   (2005). 'Complete genomic sequence of human coronavirus OC43: molecular clock analysis
474   suggests a relatively recent zoonotic coronavirus transmission event', *Journal of Virology*,
475   79/3: 1595–604. DOI: 10.1128/JVI.79.3.1595-1604.2005
476 Zhou, Z., Hui, K. P. Y., So, R. T. Y., Lv, H., Perera, R. A. P. M., Chu, D. K. W., Gelaye, E., et al.
477   (2021). 'Phenotypic and genetic characterization of MERS coronaviruses from Africa to
478   understand their zoonotic potential', *Proceedings of the National Academy of Sciences of the
479   United States of America*, 118/25: e2103984118. DOI: 10.1073/pnas.2103984118
480
481

**Supplementary Figure 1. Location of charged amino acid changes in the spike protein.** The spike protein sequences encoded by the first 300 reported genomes for the indicated SARS-CoV-2 lineages were collected, and charged amino acid changes from the original B lineage spike sequence were plotted. Charge changes were colored from dark blue (change from positive to negative charged amino acid (AA)), blue change from neutral to negative charged AA), orange (change from neutral to positive charged AA) and red (change from negative to positive charged AA). Substitutions are indicated by original AA/position in reference sequence spike/novel AA. The GenBank NC_045512 genome was used as reference. Key spike protein features of the SARS-CoV-2 spike protein are shown in the upper panel of the figure.