

1 Article

## 2 Cooccurrence of N501Y, P681R and other key mutations in 3 SARS-CoV-2 Spike

4 Carol Lee<sup>1</sup>, Shruthi Mangalaganesh<sup>2,3</sup>, Laurence O.W. Wilson<sup>1</sup>, Michael J. Kuiper<sup>4</sup>, Trevor W. Drew<sup>2,5</sup> and Seshadri  
5 S. Vasan<sup>2,6,\*</sup>

6 <sup>1</sup> Commonwealth Scientific and Industrial Research Organisation, Transformational Bioinformatics Group,  
7 North Ryde, NSW 2113, Australia; [carol.lee@csiro.au](mailto:carol.lee@csiro.au); [laurence.wilson@csiro.au](mailto:laurence.wilson@csiro.au)

8 <sup>2</sup> Commonwealth Scientific and Industrial Research Organisation, Australian Centre for Disease Prepared-  
9 ness, 5 Portarlington Road, Geelong, VIC 3220, Australia; [shruthi.mangalaganesh@csiro.au](mailto:shruthi.mangalaganesh@csiro.au); [trevor.drew@csiro.au](mailto:trevor.drew@csiro.au); [vasan.vasan@csiro.au](mailto:vasan.vasan@csiro.au)

10 <sup>3</sup> Monash University, Monash Biomedicine Discovery Institute, Clayton, VIC 3800, Australia;  
11 [sman0040@student.monash.edu](mailto:sman0040@student.monash.edu)

12 <sup>4</sup> Commonwealth Scientific and Industrial Research Organisation, Data61, Docklands, VIC 3008, Australia;  
13 [michael.kuiper@data61.csiro.au](mailto:michael.kuiper@data61.csiro.au)

14 <sup>5</sup> University of Nottingham, School of Veterinary Medicine and Science, Sutton Bonington Campus, LE12  
15 5RD, United Kingdom

16 <sup>6</sup> University of York, Department of Health Sciences, York YO10 5DD, United Kingdom,  
17 [prof.vasan@york.ac.uk](mailto:prof.vasan@york.ac.uk)

18 \* Correspondence: [vasan.vasan@csiro.au](mailto:vasan.vasan@csiro.au); Tel.: +61352275346  
19  
20

21 **Abstract:** Analysis of circa 4.2 million severe acute respiratory syndrome coronavirus 2  
22 (SARS-CoV-2) genome sequences on 'Global Initiative on Sharing All Influenza Data (GISAID)'  
23 shows the spike mutations 'N501Y' (common to Alpha, Beta, Gamma, Omicron variants) and  
24 'P681R' (central to Delta variant's spread) have cooccurred 3,678 times between 17 October 2020  
25 and 1 November 2021. In contrast, the N501Y+P681H combination is present in Alpha and Omicron  
26 variants and circa 1.1 million entries. Two-thirds of the 3,678 cooccurrences were in France,  
27 Turkey or US (East Coast), and the rest across 62 other countries. 55.5% and 4.6% of the  
28 cooccurrences were Alpha's Q.4 and Gamma's P.1.8 sub-lineages acquiring P681R; 10.7% and 3.8%  
29 were Delta's B.1.617.2 lineage and AY.33 sub-lineage acquiring N501Y; remaining 10.2% were in  
30 other variants. Despite the selective advantages individually conferred by N501Y and P681R, the  
31 N501Y+P681R combination counterintuitively didn't outcompete other variants in every instance.  
32 Although a relief to worldwide public health efforts, *in vitro* and *in vivo* studies are urgently re-  
33 quired in the absence of a strong *in silico* explanation for this phenomenon. This study demon-  
34 strates a pipeline to analyse combinations of key mutations from public domain information in a  
35 systematic manner and provide early warnings of spread.

36 **Keywords:** COVID-19; D614G; N501Y; P681R; mutations; Delta; Omicron; SARS-CoV-2; variants of  
37 concern  
38

### 39 1. Introduction

40 The severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), which causes  
41 the ongoing novel coronavirus disease 19 (COVID-19) pandemic, has a single-stranded  
42 positive-strand ribonucleic acid genome (ssRNA(+) genome) of size 26–32 kb, with high  
43 fidelity replication due to 3'-to-5' exoribonuclease 'proof-reading' mechanism [1,2]. As  
44 this virus adapts to its human host, we have seen it evolve and present quasispecies di-  
45 versity [2,3]. While most of the several thousands of mutations catalogued to date aren't  
46 substantial functional changes they have proven aetiologically useful [4,5].

47 Two years on since the start of the COVID-19 pandemic, we now have a good idea of  
48 the key mutations, especially in the Spike protein, which are punctuations in the evolu-  
49 tionary story of this virus to date. The D614G mutation reported in April 2020 resulted in  
50 the 'G-strain' with increased infectivity replacing the genomic background of this virus  
51 globally, although there was no impact on vaccine efficacy [5,6,7]. However, the N501Y  
52 mutation common to the Alpha, Beta, Gamma and Omicron variants of concern (VOC),  
53 has contributed to enhanced infection and transmission, reduced vaccine efficacy, and  
54 the ability of SARS-CoV-2 to infect new species such as wild type mice [8,9,10,11]. An-  
55 other key mutation is the P681R which alters the furin cleavage site, and has been re-  
56 sponsible for increased infectivity, transmission and global impact of the Delta variant  
57 [12,13,14].

58  
59 Our primary objective is to investigate the risk of the three aforementioned 'muta-  
60 tions of current interest' (MOCI) cooccurring naturally due to convergent evolution and  
61 resulting in a SARS-CoV-2 variant that is of greater concern than those declared to date,  
62 noting that the latest Omicron VOC has the N501Y but not the P681R mutation. Our  
63 secondary objective is to demonstrate a methodology and pipeline to analyse the spread  
64 of variants containing combinations of important mutations. We have achieved this by  
65 mining data from GISAID, the Global Initiative on Sharing All Influenza Data, which is  
66 the largest and the most comprehensive repository of SARS-CoV-2 sequences [15,16]. We  
67 have used quasispecies theory and *in silico* modelling to interpret our findings to the ex-  
68 tent possible, and recommended future research directions.

## 70 2. Results: Frequency and distribution of D614G, N501Y and P681R mutations

71  
72 The earliest record of these three MOCI coming together was in Slovenia on 17 Oc-  
73 tober 2020, however, no further observations were recorded since then in that country.  
74 Mining of GISAID data found 3,678 entries (0.1%) containing the MOCI and a majority of  
75 these were in current VOC – Alpha (61.3%), Beta (0.4%), Gamma (6.4%), Delta (21.2%),  
76 Omicron (0.0%) – and a small proportion in current or former VOI – Mu (0.3%), Kappa  
77 (0.2%) and Iota (0.02%). A small proportion (10.2%) was also observed in other variants  
78 not classified as VOC, VOI, Variant Under Monitoring (VUM), or specific sub-lineages  
79 (Figure 1).

80  
81 Figure 2 illustrates the MOCI frequency in countries with at least 50 recorded se-  
82 quences – Brazil (5.8%), Denmark (1.9%), France (22.7%), Germany (3.6%), Sweden  
83 (5.2%), Turkey (21.4%), UK (3.6%) and USA (22.1%). These eight countries represent  
84 86.4% of the cases containing the MOCI (Figure 3A, which shows a continuous timeline).  
85 Fifty-seven other countries recorded less than 50 entries for the period 17 October 2020 to  
86 15 September 2021 (Supplementary Table S1).

87  
88 Three countries – France, USA and Turkey – stand out as they each contribute to  
89 over 22% of the total instances of the MOCI cooccurring. It's unsurprising that these  
90 trends overlap with the spread of VOC in these nations (Figures 3B-3D), because there  
91 would have been more opportunities for Delta to acquire N501Y, and for Alpha or  
92 Gamma to acquire P681R (Figures 3B-3F). Such cooccurrences also appear to have hap-  
93 pened over short periods of time, for instance during March to May 2021 in France and  
94 USA with the Alpha VOC; and between June and August 2021 in Turkey with the Delta  
95 VOC (Figures 2 and 3B-3D). This could be due to founder effects, but we cannot establish  
96 this by mining GISAID data alone; it will require detailed epidemiological investigations  
97 by national public health authorities, which is beyond the scope of this work. Neverthe-  
98 less, from Figure 2, we can be reasonably sure that a number of independent events of  
99 convergent evolution (i.e. MOCI cooccurring) have taken place. It is worth investigating

100 why we do not see as many instances of Delta acquiring N501Y in France and USA, even  
101 though they had a very high number of this VOC from July 2021.  
102

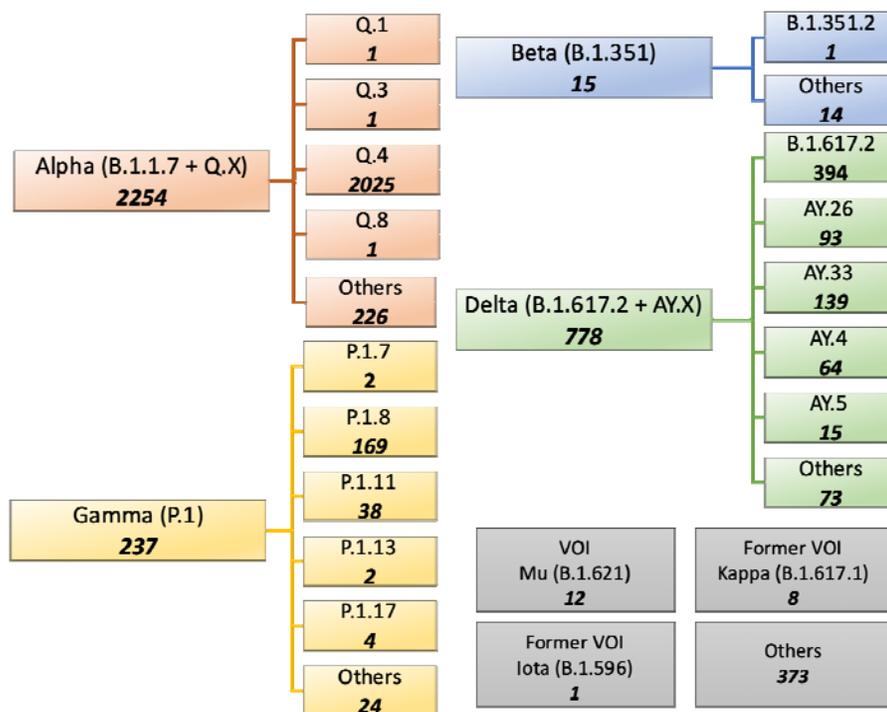


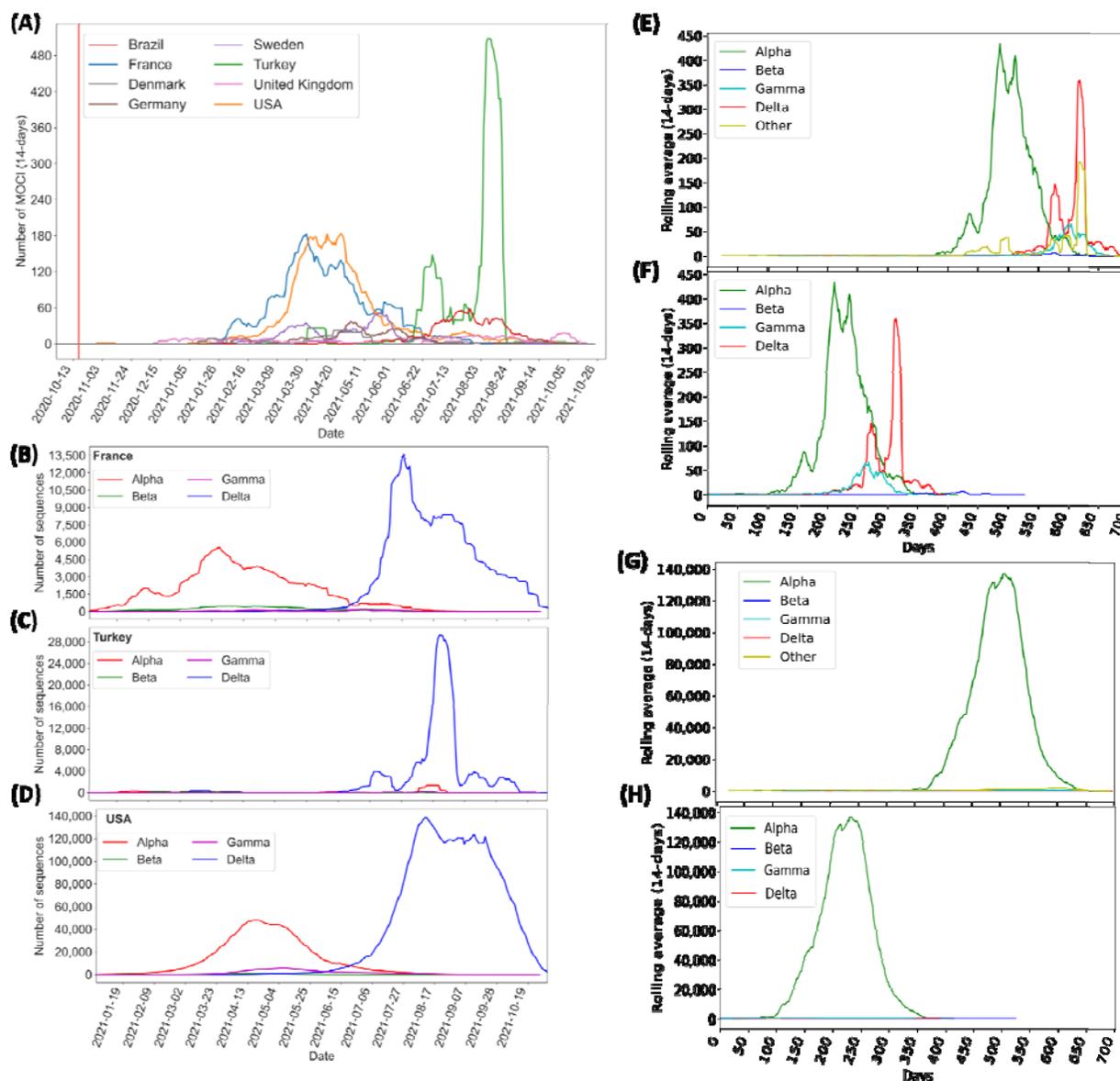
Figure 1. Summary of SARS-CoV-2 lineages and sub-lineages containing the D614G, N501Y and P681R mutations of current interest (MOCI) on GISAID. 3,678 sequences were found with our MOCI from a total of 4,177,098. Bold numbers in each box shows the number of observations for each VOC/VOI and their sub-lineages (where such information was available). 'Others' indicate samples not further classified on GISAID.

103  
104  
105  
106  
107  
108  
109



119

120



121

122

123

124

125

126

127

128

129

130

131

132

**Figure 3. Continuous representation of the timeline during which SARS-CoV-2 variants with the P681R/H, N501Y and D614G mutations were observed between October 2020 and October 2021. (A)** Number of isolates with MOCI, plotted as 14-day rolling average, in the eight countries which had at least 50 observations. The vertical red line shows the first recorded cooccurrence in Slovenia on 17 October 2020; **(B-D)** Number of VOC sequences observed in the top three countries (France, Turkey, USA respectively); **(E-F)** Number of isolates with MOCI which are also VOC, plotted as 14-day rolling average, either from the notional start of the pandemic in December 2019, or from September, May, November and October 2020 as the first reported months for Alpha, Beta, Gamma and Delta respectively. **(G-H)** Number of isolates with the P681H (rather than P681R), N501Y and D614G, plotted as 14-day rolling average, either from the notional start of the pandemic in December 2019, or from the month of first report for each VOC. ‘Other’ indicates MOCI found in variants other than these VOC.

133

134

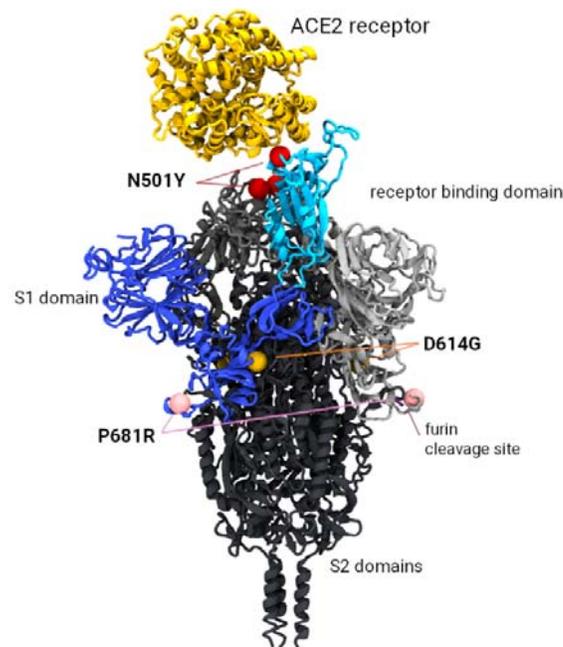
135

### 3. Discussion

136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146

### 3.1. Y501-R681 almost always in the G614 background

From **Supplementary Table S2** we see that 3,688 entries contain N501Y+P681R, just ten more than the 3,678 entries that contain D614G+N501Y+P681R; in other words, the instances of N501Y cooccurring with P681R have almost always happened on the D614G background as predicted [5,6,7]. Although these three mutations are positioned sufficiently far apart in the 3-dimensional protein structure to suggest they don't interact (**Figure 4**), further studies are required to understand whether there may be indirect functional links that could enhance viral efficiency.



147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161  
162  
163  
164  
165  
166  
167  
168  
169

**Figure 4. Model of the SARS-CoV-2 Spike trimer protein indicating the relative position of the P681R, N501Y and D614G mutations of current interest (MOCI). The N501Y mutation occurs in the receptor binding domain (highlighted in cyan) and is associated with increased binding affinity to the ACE2 receptor (shown in yellow). The P681R mutation is implicated with more efficient cleavage of the S1/S2 furin site (required prior to viral fusion to the host cell). Mutation D614G occurs at the S1/S2 interface and has also been implicated with increased replication efficiency.**

### 3.2. Predominance of Alpha Q.4 acquiring P681R and Delta acquiring N501Y

We see that the cooccurrence of the MOCI is largely due to; (a) the Alpha VOC acquiring the P681R (61.3%), and this has happened overwhelmingly in Alpha's Q.4 sub-lineage (89.9%) by definition; (b) the Delta VOC acquiring the N501Y (21.2%), and this has mainly happened in Delta's parent lineage (B.1.617.2; 50.6%) and also the sub-lineage AY.33 (17.9%); (c) the Gamma VOC acquiring the P681R (6.4%), predominantly in its P.1.8 sub-lineage (71.3%). It is worth noting that the Q.4 sub-lineage has the P681R by definition (<https://outbreak.info/situation-reports?pango=Q.4>). The Beta VOC acquiring the P681R only constituted 0.4% of the total occurrences, and this was almost entirely in the B.351.2 parent lineage (93.3%). **Supplementary Table S3** shows the results for other key Spike mutations in comparison to those obtained from GISAID. For H69del and Y145del, more samples were identified using our pipeline than on GISAID – 1 138 302 versus 1 123 439 before data filtering, and 2 563 versus 1 364 after data filtering.

170  
171 *3.3. Alpha acquiring P681R ahead of other VOC could be due to P/H681R*  
172

173 Alpha, Beta and Gamma each have N501Y, however Alpha was the first VOC to  
174 acquire P681R, both in terms of absolute (**Figure 3E**) and relative (**Figure 3F**) timelines,  
175 although Beta was reported four months before Alpha in May 2020. Alpha, especially its  
176 Q.4 sub-lineage, also contributes to most instances of the MOCI cooccurring. There is a  
177 possible biomolecular basis for this; P681H is present in Alpha, but not in the other two  
178 VOC. The Grantham scores associated with P681R (103) and P681H (77) are comparable  
179 and the H681R substitution is thus a conservative change (Grantham distance 29). There  
180 is only one way to get to Histidine (H) or Arginine (R) from Proline (P); and for H and R  
181 there is only one way to get to each other; and we see no structural reason from our *in*  
182 *silico* model as to why one should be preferred to the other. Thus, we could infer that  
183 some of the instances of Alpha acquiring P681R could have been due to a substitution of  
184 H with R, which is the signature of the Q.4 sub-lineage. Spike's 681st position is the fifth  
185 substrate sequence for cleavage recognition; both furin and the transmembrane serine  
186 protease 2 (TMPRSS2) cleave the Spike at 685/686 position with H (and possibly R) en-  
187 hancing this process [7,17,18]. It is worth noting that to penetrate host cells, the  
188 SARS-CoV-2 Spike protein must be cut twice by host proteins. In the SARS-CoV-1  
189 (SARS), both incisions occur after the virus has locked on to a cell. But with SARS-CoV-2,  
190 the presence of the furin cleavage site enables host enzymes like furin to make the first  
191 cut as newly formed viral particles emerge from an infected cell. These pre-activated viral  
192 particles can then go on to infect cells more efficiently compared to particles requiring  
193 two cuts. Thus, the P681R increases the susceptibility of the furin cleavage site in Delta  
194 VOC, and allows the exposure of the Spike's S2 subunit for better cell integration.  
195

196 *3.4. Cooccurrences have been reported predominantly in eight countries*

197 Cooccurrences of the MOCI were observed with Alpha, Beta, Gamma and Delta  
198 VOC in 41, 8, 11 and 47 countries respectively, indicating convergent evolution, espe-  
199 cially as much of the world was under lockdown during 2020-2021. Curiously, two thirds  
200 of the observations have been reported from just France, Turkey and USA; 86.3% from  
201 these three countries plus Brazil, Denmark, Sweden, UK and Germany (**Supplementary**  
202 **Table S1; Supplementary Figure S4**). Several instances involve proximal regions,  
203 suggesting multiple founder effects (**Figures 2 and 5**). For example, in Denmark and  
204 Sweden, variants with the MOCI were concentrated in highly populous Hovedstaden  
205 and Svealand regions (Sörmland, Stockholm, Uppsala and Västmanland) [19]; the  
206 MOCI frequencies were also correlated with the most common lineage in respective  
207 countries and regions, suggesting community transmission, although these MOCI did  
208 not go on to become the dominant variant in those locations. This is in line with most  
209 within-host variants getting lost during transmission, and only a few founding infections  
210 being maintained in a given population [20].

211 Cooccurrence of the MOCI in the UK occurred mostly with the Alpha VOC's B.1.1.7  
212 and Q.4 from late 2020 to May 2021; followed by the second wave of the Delta VOC since  
213 April 2021, especially in AY.4 that had the highest (58%) prevalence (**Figure 2**). The for-  
214 mer period had ten times the death rate (~1200 per week between September 2020 and  
215 March 2021) compared to the latter (~120 per week; probably due to lockdowns and 30%  
216 of the population already double-vaccinated). In the UK which has a very high rate of  
217 sequencing, the overall numbers of cooccurring MOCI were still low (133 records, 88% in  
218 England). Of these, there were respectively 52, 39 and 18 instances of the AY.4 acquiring  
219 N501Y, and the B.1.1.7 and Q.4 acquiring P681R [21,22]. It is possible that the lockdown  
220 protocols and vaccination coverage may have influenced the observed viral transmission  
221 dynamics [23]. In the Americas, most of 1,078 cooccurrences of the MOCI were detected

222 either in the US East Coast (673 entries), or in the Sul (Santa Catarina) and Sudeste (São  
223 Paulo, Rio de Janeiro) regions in the South and Southeast of Brazil (284 entries). These  
224 trends could be due to founder effects and lack of travel restrictions from early 2021, es-  
225 pecially in the United States [24] (Figures 2 and 5). In Brazil, the cooccurrences of the  
226 MOCI also corresponded with the prevalence of P.1 until August 2021 after which the  
227 AY.99.2 became the dominant variant (Supplementary Figure S4). Unfortunately, re-  
228 gional and city level information was not available for Turkey on GISAID, preventing  
229 further analysis and insights and highlighting the importance of metadata [16].

230

231

232

233

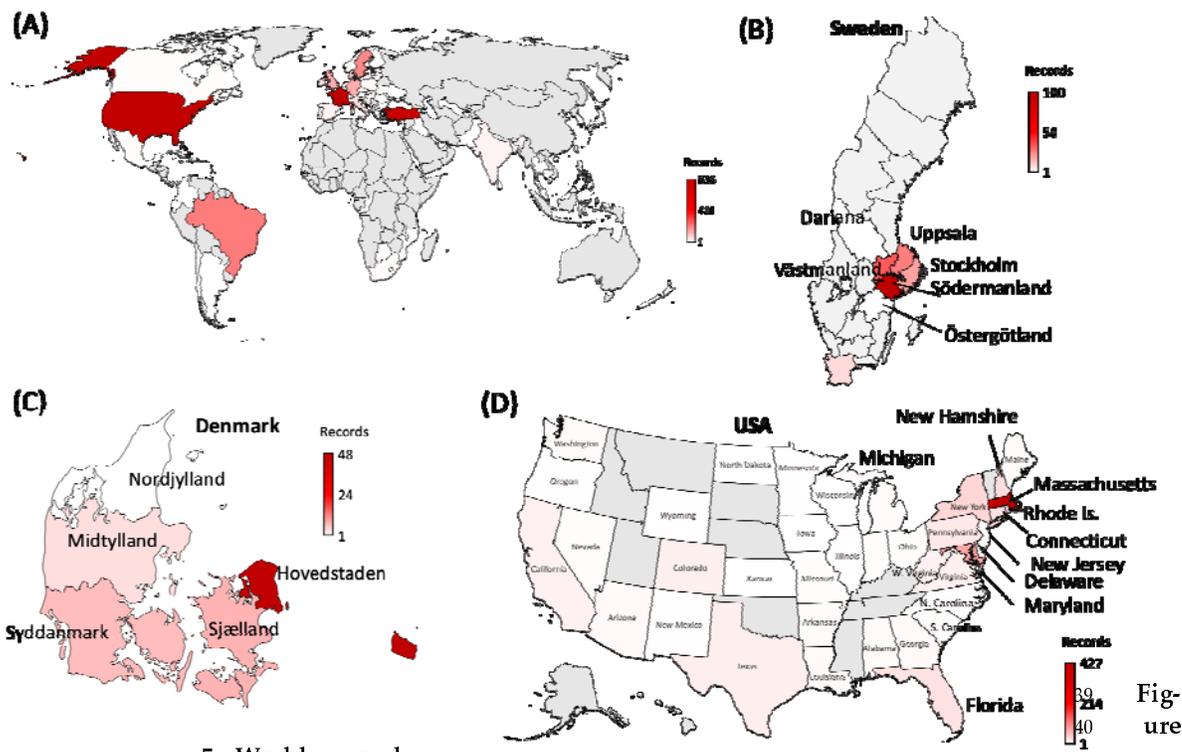
234

235

236

237

238



241

242 **5. World map de-**  
243 **depicting the 65 countries where SARS-CoV-2 variants containing the cooccurring MOCI**  
244 **were recorded between October 2020 and October 2021. (A) world map, with zoomed in**  
245 **view of the regions observed with MOCI shown for (B) Sweden (C) Denmark and (D)**  
246 **USA. Countries are color-coded based on the number of records of MOCI observed, as**  
247 **shown in the legend, with greyed out regions indicating no observations. For some**  
248 **countries such as Turkey, region/city-level information was not available, but EPI\_ISL**  
249 **numbers can be provided upon request.**

249

### 250 3.5. Y501-R681 doesn't outcompete other variants

251

252 From our analysis, SARS-CoV-2 isolates containing Y501-R681 in Spike did not  
253 outcompete other variants in every single instance we have examined (Figures 2 and 3),  
254 which is both counterintuitive and a relief to worldwide public health efforts. Although  
255 individually these two mutations have been reported to confer selective advantage  
256 [8,9,10,11,12,13,14], their combination apparently has not and we have not found any  
257 compelling biomolecular reason for this (Figure 4). It would be desirable to conduct *in*  
258 *vitro* and *in vivo* studies to gain a better understanding, either using infectious clones and  
259 reverse genetics [25,26,27,28], or using naturally occurring comparable isolates with and  
260 without these mutations [see for instance 7]

261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
  
305  
306  
307  
308  
309  
310  
311  
312

### 3.6. Implications for the Omicron VOC

Alpha and Omicron each have the Y501-H681 combination which is present in 1,085,434 entries out of the 4.2 million we have examined in this study (i.e. over 25%); this contrasts with the Y501-R681 of which we have only found 3,688 observations. Just as the Q.4 sub-lineage with Y501-R681 emerged when the Alpha VOC started to spread, it is possible that a new sub-lineage of Omicron with Y501-R681 could emerge given the rapid rise of this latest VOC [29,30,31]. In this context, it is important to highlight the problem of sequencing artefact, for example 65 high coverage sequences of Omicron with P681R recently appeared on GISAID but were subsequently corrected by the submitting laboratory following repeat libraries with fewer cycles and other improvements [32]. The spread of VOC in immunocompromised patients could lead to mutations with a selective advantage for antibody escape and/or transmissibility (e.g. N501Y and P681H/R), in addition to a number of deletions in the N terminal domain [33], and this aspect needs further investigation given Omicron's large number of mutations (many of them yet to be studied in depth). The adaptation and evolution of this virus in new hosts, for instance mice and rats, also pose additional risks such as additional reservoirs and reinfection of humans [11]. Investigations into the dynamics affecting the SARS-CoV-2 pandemic in each of the eight countries are also warranted, in the wider context of travel restrictions, population demographics and dynamics [34,35]. Lack of patient-deidentified metadata in a consistent format further complicates meaningful analysis as emphasized before [16].

### 3.7. SARS-CoV-2 evolution and selective pressures on quasispecies

Any virus which is new to a host will undergo a period of adaptation, as numerous selective pressures come to bear. RNA viruses, such as coronaviruses, are error-prone in their replication and exist as a cloud of variants, with the dominant clade representing the majority population, but with a number of diverse variants also represented within the population, known as quasispecies [2,3,36,37]. Among viruses which are established in their host, the dominant clade is generally relatively 'fit' in its ability to replicate and out-competes all the other clades. However, when a virus is new to a host, its relative fitness may be quite low, so, in any particular environment, different clades can have similar relative fitness. This means subtle differences might give a particular clade a small advantage enabling it to predominate as the population continues its host adaptation. The environment may comprise of host factors, such as host genetics, tissue tropism, age, immune status or competence, presence of other infections, as well as external factors, such as temperature and humidity. With SARS-CoV-2, the extraordinary public health measures in some regions of the world are likely exerting their own selective pressures on its evolution and creating numerous local environments. This phenomenon, called the 'survival of the flattest' [3] can lead to the emergence of multiple clades dominating under different environments. There is also an increasing appreciation that, within the host, different clades may predominate in different tissues at the same time and that there may be interactions among the diverse quasispecies that facilitate the infection [38].

The pathway of adaptation is not always towards severe disease, since this may not lead to maximum replication and transmission. Indeed, many viruses with a long history of acute infection in an established host species, the disease may be mild or inapparent, reflecting a relationship more towards stasis, where the impact on the host does not compromise the chances of the progeny virus being transmitted to a new host. For many coronaviruses, we have seen an evolution towards milder infection, albeit with greater transmission rates than has been seen with SARS-CoV-2 to date. Certainly, some level of reduction in the efficacy of vaccine-induced antibody in neutralising the Delta

313 and Omicron variants has been demonstrated, but the other feature of these viruses,  
314 which is less often discussed, is the ability of these VOC to replicate by means of syncytial  
315 formation, rather than the simple apoptotic cycle seen with earlier variants [39]. This  
316 may confer an important immune evasive aspect to replication, which might offer a dis-  
317 tinct advantage, in immune avoidance. It might also explain why these clades are more  
318 productive and less pathogenic since systemic infection is not always involved.

### 319 320 *3.8. Error catastrophe, Muller's ratchet and implications for SARS-CoV-2*

321  
322 A phenomenon related to quasispecies evolution is that of error catastrophe, which  
323 may occur if the mutational rate exceeds the ability of the fitness landscape to  
324 accommodate the resultant evolving clades [40]. The key triggers for this are not fully  
325 understood, but it is thought that the application of strict biosecurity measures may play  
326 a part by restricting the availability of new hosts. In such cases, the quasispecies becomes  
327 increasingly attenuated, with mutations accumulating via a mechanism called 'Muller's  
328 ratchet', to the point where the epidemic may die out [41]. This has been seen in many  
329 prior outbreaks of Ebola, where the infection became increasingly less pathogenic and  
330 less transmissible. While the current environment of SARS-CoV-2 is complex and highly  
331 variable, the selective pressure on a virus is towards higher transmissibility and milder  
332 disease in all environments, and the current systems of control preferentially select for  
333 such. Once a virus develops either/both of these traits, it is very unusual for the reverse  
334 to occur because viruses with higher transmissibility out-compete those with lesser  
335 transmissibility and those which are milder disable the host less, so they continue to  
336 travel and mix with susceptible hosts. It remains to be seen if this is being observed with  
337 the Omicron VOC as early reports indicate higher transmissibility and less severe  
338 disease.

339  
340 In comparing variants of SARS-CoV-2, it is tempting to focus on specific mutations  
341 and their location – also perhaps with overt attention being paid to the Spike protein,  
342 especially its receptor binding domain – and not consider that the observed changes  
343 could evolve independently of each other. The secondary and tertiary folding of proteins,  
344 as well as post translational modifications, can have a profound effect on function, such  
345 that a mutation at one site might require a complimentary change at another site in order  
346 to confer advantage, or be wiped out by another, seemingly unrelated change. For the  
347 same reason, convergent evolution, might also occur, exhibited by the co-existence of  
348 identical co-mutations in virus clades of different lineage [42]. Another phenomenon  
349 known to occur in many coronaviruses is recombination, where partial genomic  
350 exchange may occur between two coronaviruses when simultaneously infecting the same  
351 host. The extent to which this may have occurred in SARS-CoV-2 is the likely subject of  
352 future analysis [43,44].

## 353 354 **4. Methods**

### 355 356 *4.1. Bioinformatics analysis*

357  
358 We mined GISAID data on 3 November 2021 containing 4,835,087 entries, 99.9%  
359 (4,831,865) sequenced from people who tested COVID-19 positive. Of these, 4,748,462  
360 contained over 29 thousand nucleotides and were deemed 'complete'. GISAID has also  
361 defined a subset of these 4.8 million sequences, containing 3,519,085 entries, as 'high  
362 coverage'. This is because 99% of their bases are defined; there is no unverified deletion  
363 or insertion; and unique mutations not seen in other sequences constitute fewer than

364 0.05% in these entries. Consequently, some of our analysis of the data downloaded on 1  
365 November 2021 (4,766,139 sequences) required manual curation whenever we encoun-  
366 tered incorrect dates, or low coverage sequencing that could skew our results. **Supple-**  
367 **mentary Table S5** provides GISAID's break down of these 4.8 million entries by variants  
368 containing the D614G (98.8%), N501Y (27.01%) and P681R (48.2%) mutations; by Alpha  
369 (23.9%), Beta (0.8%), Gamma (2.4%), Delta (47.8%), Omicron (0.0%) VOC; and by Lambda  
370 (0.2%) and Mu (0.3%) variants of interest (VOI).

371  
372 Human-origin SARS-CoV-2 sequences over 29 kb length were aligned to the nomi-  
373 nal reference (EPI\_ISL\_402123) using an in-house alignment pipeline that generated a file  
374 in the 'variant call format' (VCF) containing all the mutations as of 1 November 2021  
375 [1,2]. In the VCF containing 4,294,469 genome sequences, we searched for entries with  
376 MOCI, for example, the Spike N501Y was searched at nucleotide position 23403 for an A  
377 to G transition. The EPI\_ISL sample IDs containing our MOCI were merged with their  
378 respective GISAID metadata (location, date of collection, lineage) to create an annotated  
379 database. **Supplementary Table S5** shows the split up of these 4.3 million sequences by  
380 variants containing the D614G (98.8%), N501Y (29.2%) and P681R (44.2%) mutations; and  
381 by Alpha (26.3%), Beta (0.9%), Gamma (2.2%), Delta (43.8%), Omicron (0.0%), Lambda  
382 (0.2%) and Mu (0.3%) VOC/VOI.

383  
384 Using a custom Python code, we identified the presence of the D614G, N501Y and  
385 P681R mutations, and combinations thereof (c.f. 'Data Availability Statement' for our  
386 code). Our final dataset contained 4,177,098 entries after filtering out GISAID records  
387 with incomplete dates, i.e. retaining those with a 'YYYY-MM-DD' format. We ascertained  
388 whether these isolates are one of the VOC (and if yes, which sub-lineage), their location  
389 (country, state, and city depending on information available), and sample collection date.  
390 Using the latter, we calculated rolling averages over 14-days, as this window has been  
391 shown from previous experience to be an optimal period to reduce background noise,  
392 especially as the data is discrete and highly variable in size across countries [11]. How-  
393 ever, unlike our previous study in which we were interested in macroscopic trends, here  
394 we have looked at every instance of key mutations cooccurring; therefore, we did not use  
395 a threshold or filter based on minimal sample size. Our analysis presents the spread of  
396 these mutations from the notional start of the SARS-CoV-2 pandemic (31 December 2019)  
397 to the end of our observation period (01 November 2021) with appropriate classifications  
398 (e.g. Alpha, Beta, Gamma, Delta).

#### 400 4.2. Biomolecular modelling

401 Fully glycosylated *in silico* models of the SARS-CoV-2 Spike protein were con-  
402 structed based on the '6VSB' protein databank (PDB) structure [45], minus the  
403 transmembrane domain (residues 1161-1272) for additional computational efficiency as  
404 described previously [7]. Models were simulated in aqueous solution (TIP3, water,  
405 0.15M ions, NVT ensemble, 310K) using NAMD 2.14 software [46]. Models were visual-  
406 ized with the visual molecular dynamics (VMD) program for large biomolecular systems  
407 which uses 3-D graphics and built-in scripting to assess the spatial arrangement of mu-  
408 tations and any complementary or inhibitory interactions [47].

## 410 5. Conclusions and limitations

411  
412 SARS-CoV-2 is the most sequenced virus in the world, however, there is an inherent  
413 bias introduced by the highly variable sequencing-to-COVID-19 positivity ratio (S:P)  
414 across countries (e.g., the UK has one of the highest P:S ratios in the world), and across  
415 time (S:P has generally gone up in all countries, e.g. India). Thus, the GISAID data is

416 likely to contain non-random samples with a skew in favour of sequences associated with  
417 epidemiologically consequent cases, albeit this is difficult to decipher as most sequences  
418 have little or no meaningful patient-deidentified information. The lockdowns and vac-  
419 cination coverage associated with this viral pandemic have also been unprecedented in  
420 human history and spatio-temporally variable, further complicating our analysis. Not-  
421 withstanding these limitations, which have led to the proliferation of certain clades in  
422 certain environments, we are still able to make broad conclusions thanks to the very large  
423 number of sequences available on GISAID. We are also able to demonstrate a process and  
424 pipeline to analyse mutations of structural, functional and epidemiological consequence  
425 from public domain information in a systematic manner and provide early warnings of  
426 certain combinations starting to spread. Though it is unclear which combinations of  
427 mutations provide the best selective advantage, there is concern that mutations previ-  
428 ously observed in other dominant lineages may arise spontaneously in contemporary  
429 strains further increasing their potential for infectivity and adverse clinical outcomes. In  
430 this case-study we analysed three key mutations (viz. D614G, N501Y and P681R)  
431 cooccurring but found no evidence of this combination spreading. Although 3,678 se-  
432 quences were found on GISAID between 17 October 2020 to 1 November 2021, mainly in  
433 France, Turkey and USA, the Y501-R681 combination hasn't outcompeted other variants  
434 and this warrants further *in silico*, *in vitro* and *in vivo* investigations. With the latest VOC  
435 Omicron containing a large number of mutations, many yet to be studied in-depth, our  
436 methodology will be very useful to understand whether certain combinations of muta-  
437 tions are more transmissible. If GISAID entries are double-checked for sequencing arte-  
438 fact before submission and strengthened with patient-deidentified metadata, then this  
439 approach could enable early epidemiological intelligence, for instance on case severity,  
440 mortality, and factors such as age, gender, race and co-morbidities that could increase the  
441 infection risk.

442 **Supplementary Materials:** The following supporting information can be downloaded at:  
443 [www.mdpi.com/xxx/s1](http://www.mdpi.com/xxx/s1), Table S1: Raw data for countries with cooccurring P681R, N501Y and  
444 D614G mutations. [www.mdpi.com/xxx/s2](http://www.mdpi.com/xxx/s2), Table S2: Observation of key mutations and their com-  
445 binations on GISAID. [www.mdpi.com/xxx/s3](http://www.mdpi.com/xxx/s3), Table S3: Observation of other key SARS-CoV-2  
446 Spike mutations. [www.mdpi.com/xxx/s4](http://www.mdpi.com/xxx/s4), Figure S4: Frequency of the P681R, N501Y and D614G  
447 mutations cooccurring in the eight countries with more than 50 observations between October 2020  
448 and October 2021. [www.mdpi.com/xxx/s5](http://www.mdpi.com/xxx/s5), Table S5: GISAID data on key SARS-CoV-2 mutations  
449 classified by VOC/VOI.

450 **Author Contributions:** Conceptualization and methodology, S.S.V.; bioinformatics, C.L. and  
451 L.O.W.W.; biomolecular modelling and interpretation, M.J.K. and T.W.D.; data analysis and orig-  
452 inal draft preparation, S.M., C.L. and S.S.V.; funding acquisition, S.S.V. and L.O.W.W.; manuscript  
453 review and editing, all authors.

454 **Funding:** This work was supported by funding (Principal Investigator: S.S.V.) from the CSIRO  
455 Future Science Platforms, National Health and Medical Research Council (MRF2009092), and  
456 United States Food and Drug Administration (FDA) Medical Countermeasures Initiative contract  
457 (75F40121C00144); and by funding (Principal Investigator: L.O.W.W.) from the Australian Acad-  
458 emy of Science and the Australian Department of Industry, Science, Energy and Resources. The  
459 article reflects the views of the authors and does not represent the views or policies of the funding  
460 agencies including the FDA.

461 **Institutional Review Board and Informed Consent Statements:** Not applicable for the following  
462 reasons. This study does not involve any human samples, identifiable human data or patient re-  
463 cords. It only involves bioinformatics analysis of novel coronavirus genome sequences openly ac-  
464 cessible from the public repository 'GISAID', and these sequences do not contain any metadata that  
465 can be identified with any individual. This includes sample collection dates and anonymous sam-  
466 ple IDs that are openly accessible. Our study does not report any exact age or photographs. Clinical  
467 trial registration is not applicable as this is neither an interventional study nor an observational  
468 study involving humans.

469 **Data Availability Statement:** This paper only uses publicly accessible data that can be downloaded  
470 from GISAID, the world's largest repository of SARS-CoV-2 genome sequences. The lists of  
471 GISAID's EPI\_ISL\_ sample IDs and data used by our analysis can be requested from the corre-  
472 sponding author. Our Python code is available at  
473 <https://github.com/Carol-Lee-gh/Covid-Mutation-Pipeline.git>.

474 **Acknowledgments:** We are grateful for support from our CSIRO colleagues at the Australian  
475 Centre for Disease Preparedness (<https://www.grid.ac/institutes/grid.413322.5>), especially Kim  
476 Blasdell, Simran Chahal, Alexander McAuley and Nagendrakumar Singanallur, and the Australian  
477 e-Health Research Centre, especially Denis Bauer, David Hansen, Yatish Jain, Brendan Hosking  
478 and Aidan Tay. We thank Professor Adriana Heguy of New York University Langone Health for  
479 discussions on sequences of Omicron with P681R.

480 **Conflicts of Interest:** The authors declare no conflict of interest.

## 481 References

- 482
- 483 1. Zhou, P.; Yang, X.L.; Wang, X.G.; et al. A pneumonia outbreak associated with a new coronavirus of probable bat origin. *Nature* **2020**, *579*(7798), 270-273, <https://doi.org/10.1038/s41586-020-2012-7>.
  - 484 2. Bauer, D.C.; Tay, A.P.; Wilson, L.O.W.; et al. Supporting pandemic response using genomics and bioinformatics: A case  
485 study on the emergent SARS-CoV-2 outbreak. *Transboundary and Emerging Diseases* **2020**, *67*(4), 1453-1462,  
486 <https://doi.org/10.1111/tbed.13588>.
  - 487 3. Wilke, C.O.; Wang, J.L.; Ofria, C.; Lenski, R.E.; Adami, C. Evolution of digital organisms at high mutation rates leads to  
488 survival of the flattest. *Nature* **2020**, *412*(6844), 331-333, <https://doi.org/10.1038/35085569>.
  - 489 4. Grubaugh, N.D.; Petrone, M.E.; Holmes, E.C. We shouldn't worry when a virus mutates during disease outbreaks. *Nature*  
490 *Microbiology* **2020**, *5*(4), 529-530, <https://doi.org/10.1038/s41564-020-0690-4>.
  - 491 5. Callaway, E. The coronavirus is mutating — does it matter? *Nature* **2020**, *585*(7824), 174-177,  
492 <https://doi.org/10.1038/d41586-020-02544-6>.
  - 493 6. Korber, B.; Fischer, W.M.; Gnanakaran, S.; et al. Tracking changes in SARS-CoV-2 Spike: evidence that D614G increases  
494 infectivity of the COVID-19 virus. *Cell* **2020**, *182*(4), 812-827.e19, <https://doi.org/10.1016/j.cell.2020.06.043>.
  - 495 7. McAuley, A.J.; Kuiper, M.J.; Durr, P.A.; et al. Experimental and *in silico* evidence suggests vaccines are unlikely to be af-  
496 fected by D614G mutation in SARS-CoV-2 Spike protein. *npj Vaccines* **2020**, *5*, 96,  
497 <https://doi.org/10.1038/s41541-020-00246-8>.
  - 498 8. Liu, Y.; Liu, J.; Plante, K.S.; et al. The N501Y Spike substitution enhances SARS-CoV-2 infection and transmission. *Nature*  
499 **2021**, <https://doi.org/10.1038/s41586-021-04245-0>.
  - 500 9. Abdool Karim, S.S.; de Oliveira, T. New SARS-CoV-2 variants — clinical, public health, and vaccine implications. *New*  
501 *England Journal of Medicine* **2021**, *384*(19), 1866-1868, <http://doi.org/10.1056/NEJMc2100362>.
  - 502 10. Riddell, S.; Goldie, S.; McAuley, A.J.; et al. Live virus neutralisation of the 501Y.V1 and 501Y.V2 SARS-CoV-2 variants fol-  
503 lowing INO-4800 vaccination of ferrets. *Frontiers in Immunology* **2021**, *12*, 694857,  
504 <https://doi.org/10.3389/fimmu.2021.694857>.
  - 505 11. Kuiper, M.J.; Wilson, L.O.W.; Mangalaganesh, S.; Lee, C.; Reti, D.; Vasan, S.S. But Mouse, you are not alone: On some se-  
506 vere acute respiratory syndrome coronavirus 2 variants infecting mice. *ILAR Journal* **2021**,  
507 <http://dx.doi.org/10.1093/ilar/ilab031>.
  - 508 12. Callaway, E. The mutation that helps Delta spread like wildfire. *Nature* **2021**, *596*(7873), 472-473,  
509 <https://doi.org/10.1038/d41586-021-02275-2>.
  - 510 13. Liu, Y.; Liu, J.; Johnson, B.A.; et al. Delta Spike P681R mutation enhances SARS-CoV-2 fitness over Alpha variant. *bioRxiv*  
511 **2021**, <https://doi.org/10.1101/2021.08.12.456173>.
  - 512 14. Peacock, T.P.; Sheppard, C.M.; Brown, J.C.; et al. The SARS-CoV-2 variants associated with infections in India, B.1.617,  
513 show enhanced Spike cleavage by furin. *bioRxiv* **2021**, <https://doi.org/10.1101/2021.05.28.446163>.
  - 514 15. Elbe, S.; Buckland-Merrett, G. Data, disease and diplomacy: GISAID's innovative contribution to global health. *Global*  
515 *Challenges* **2017**, *1*(1), 33-46, <https://doi.org/10.1002/gch2.1018>.
  - 516 16. Bauer, D.C.; Metke-Jimenez, A.; Maurer-Stroh, S.; et al. Interoperable medical data: The missing link for understanding  
517 COVID-19. *Transboundary and Emerging Diseases* **2021**, *68*(4), 1753-1760, <https://doi.org/10.1111/tbed.13892>.
  - 518 17. Remacle, A.G.; Shiryayev, S.A.; Oh, E.S.; et al. Substrate cleavage analysis of furin and related proprotein convertases: A  
519 comparative study. *The Journal of Biological Chemistry* **2008**, *283*(30), 20897-20906, <https://doi.org/10.1074/jbc.M803762200>.
  - 520 18. Mahoney, M.; Damalanka, V.C.; Tartell, M.A.; et al. A Novel Class of TMPRSS2 inhibitors potently block SARS-CoV-2 and  
521 MERS-CoV viral entry and protect human epithelial Lung Cells. *bioRxiv* **2021**, <https://doi.org/10.1101/2021.05.06.442935>.
  - 522 19. Demographics | Statista, 2016-2020. Available online: <https://www.statista.com/markets/411/topic/446/demographics/>  
523 (accessed on 22 December 2021).
  - 524

- 525 20. Lythgoe, K.A.; Hall, M.; Ferretti, L.; et al. SARS-CoV-2 within-host diversity and transmission. *Science* **2021**, 372(6539),  
526 eabg0821, <https://doi.org/10.1126/science.abg0821>.
- 527 21. Latif, A.A.; Mullen, J.L.; Alkuzweny, M.; et al. outbreak.info. Available online:  
528 <https://outbreak.info/location-reports?loc=GBR&selected=Delta&selected=Alpha> (accessed on 22 December 2021).
- 529 22. Ritchie, H.; Mathieu, E.; Rodés-Guirao, L.; et al. Coronavirus pandemic (COVID-19), 2020. OurWorldInData. Available  
530 online: <https://ourworldindata.org/coronavirus> (accessed on 22 December 2021).
- 531 23. Timeline of UK Government Coronavirus Lockdowns, 2021. The Institute for Government. Available online:  
532 <https://www.instituteforgovernment.org.uk/charts/uk-government-coronavirus-lockdowns> (accessed on 22 December  
533 2021).
- 534 24. Travel restrictions issued by states in response to the coronavirus (COVID-19) pandemic, 2020-2021. Ballotpedia. Available  
535 online:  
536 [https://ballotpedia.org/Travel\\_restrictions\\_issued\\_by\\_states\\_in\\_response\\_to\\_the\\_coronavirus\\_\(COVID-19\)\\_pandemic,\\_20\\_20-2021](https://ballotpedia.org/Travel_restrictions_issued_by_states_in_response_to_the_coronavirus_(COVID-19)_pandemic,_20_20-2021)  
537 (accessed on 22 December 2021).
- 538 25. Thao, T.T.N.; Labroussaa, F.; Ebert, N.; et al. Rapid reconstruction of SARS-CoV-2 using a synthetic genomics platform.  
539 *Nature* **2020**, 582(7813), 561-565, <https://doi.org/10.1038/s41586-020-2294-9>.
- 540 26. Xie, X.; Muruato, A.; Lokugamage, K.G.; et al. An infectious cDNA clone of SARS-CoV-2. *Cell Host & Microbe* **2020**, 27(5),  
541 841-848, <https://doi.org/10.1016/j.chom.2020.04.004>.
- 542 27. Hou, Y.J.; Okuda, K.; Edwards, C.E.; et al. SARS-CoV-2 reverse genetics reveals a variable infection gradient in the respira-  
543 tory tract. *Cell* **2020**, 182(2), 429-446, <https://doi.org/10.1016/j.cell.2020.05.042>.
- 544 28. Xie, X.; Lokugamage, K.G.; Zhang, X.; et al. Engineering SARS-CoV-2 using a reverse genetic system. *Nature Protocols* **2021**,  
545 16(3), 1761-1784, <https://doi.org/10.1038/s41596-021-00491-8>.
- 546 29. Ferguson, N.; Ghani, A.; Cori, A.; Hogan, A.; Hinsley, W.; Volz, E. Report 49 - Growth, population distribution and im-  
547 mune escape of Omicron in England. Imperial College London. Available online:  
548 <https://www.imperial.ac.uk/mrc-global-infectious-disease-analysis/covid-19/report-49-Omicron/> (accessed on 22 Decem-  
549 ber 2021).
- 550 30. Cele, S.; Jackson, L.; Khoury, D.S.; et al. SARS-CoV-2 Omicron has extensive but incomplete escape of Pfizer BNT162b2  
551 elicited neutralization and requires ace2 for infection. *medRxiv* **2021**, <https://doi.org/10.1101/2021.12.08.21267417>.
- 552 31. Pulliam, J.R.C.; Schalkwyk, C.V.; Govender, N.; et al. Increased risk of SARS-CoV-2 reinfection associated with emergence  
553 of the Omicron variant in South Africa. *medRxiv* **2021**, <https://doi.org/10.1101/2021.11.11.21266068>.
- 554 32. Heguy, A. (NYU Langone Health, New York, NY, USA). Personal communication to S.S. Vasan, 30 December 2021.
- 555 33. Corey, L.; Beyrer, C.; Cohen, M.S.; et al. SARS-CoV-2 variants in patients with immunosuppression. *New England Journal of*  
556 *Medicine* **2021**, 385(6), 562-566, <https://doi.org/10.1056/NEJMs2104756>.
- 557 34. Grépin, K.A.; Ho, T.L.; Liu, Z.; Marion, S.; Piper, J.; Worsnop, C.Z.; Lee, L. Evidence of the effectiveness of travel-related  
558 measures during the early phase of the COVID-19 pandemic: A rapid systematic review. *BMJ Global Health* **2021**, 6(3),  
559 e004537, <https://doi.org/10.1136/bmjgh-2020-004537>.
- 560 35. Liebig, J.; Najeebullah, K.; Jurdak, R.; Shoghri, A.E.; Paini, D. Should international borders re-open? The impact of travel  
561 restrictions on COVID-19 importation risk. *BMC Public Health* **2021**, 21(1), 1573, <https://doi.org/10.1186/s12889-021-11616-9>.
- 562 36. Domingo, E.; Perales, C. Viral quasispecies. *PLoS Genetics* **2019**, 15(10), e1008271,  
563 <https://doi.org/10.1371/journal.pgen.1008271>.
- 564 37. Drew, T.W. The emergence and evolution of swine viral diseases: To what extent have husbandry systems and global trade  
565 contributed to their distribution and diversity. *Revue Scientifique et Technique de l'OIE* **2011**, 30(1), 95-106,  
566 <https://doi.org/10.20506/rst.30.1.2020>.
- 567 38. Sun, F.; Wang, X.; Tan, S.; et al. SARS-CoV-2 quasispecies provides an advantage mutation pool for the epidemic variants.  
568 *Microbiology Spectrum* **2021**, 9(1), e00261-21, <https://doi.org/10.1128/Spectrum.00261-21>.
- 569 39. Li, L.; Li, Q.; Shi, Y. Syncytia formation during SARS-CoV-2 lung infection: a disastrous unity to eliminate lymphocytes.  
570 *Cell Death & Differentiation* **2021**, 28, 2019-2021, <https://doi.org/10.1038/s41418-021-00795-y>
- 571 40. Tejero, H.; Marín, A.; Montero, F. The relationship between the error catastrophe, survival of the flattest, and natural se-  
572 lection. *BMC Evolutionary Biology* **2011**, 11(1), 2, <https://doi.org/10.1186/1471-2148-11-2>.
- 573 41. de Alcañiz, J.G.G.; López-Rodas, V.; Costas, E. Sword of Damocles or choosing well: Population genetics sheds light into  
574 the future of the COVID-19 pandemic and SARS-CoV-2 new mutant strains. *medRxiv* **2021**,  
575 <https://doi.org/10.1101/2021.01.16.21249924>.
- 576 42. Brüßow, H. COVID-19: emergence and mutational diversification of SARS-CoV-2. *Microbial Biotechnology* **2021**, 14(3),  
577 756-768, <https://doi.org/10.1111/1751-7915.13800>.
- 578 43. Haddad, D.; John, S.E.; Mohammad, A.; et al. SARS-CoV-2: Possible recombination and emergence of potentially more  
579 virulent strains. *PLoS One* **2021**, <https://doi.org/10.1371/journal.pone.0251368>.
- 580 44. Pollett, S.; Conte, M.A.; Sanborn, M.; et al. A comparative recombination analysis of human coronaviruses and implications  
581 for the SARS-CoV-2 pandemic. *Scientific Reports* **2021**, 11, 17365, <https://doi.org/10.1038/s41598-021-96626-8>.
- 582 45. Wrapp, D.; Wang, N.; Corbett, K.S.; et al. Cryo-EM structure of the 2019-nCoV Spike in the prefusion conformation. *Science*  
583 **2020**, 367(6483), 1260-1263, <https://doi.org/10.1126/science.abb2507>.

- 584 46. Phillips, J.C.; Braun, R.; Wang, W.; et al. Scalable molecular dynamics with NAMD. *Journal of Computational Chemistry* **2005**,  
 585 26(16), 1781-1802, <https://doi.org/10.1002/jcc.20289>.  
 586 47. Humphrey, W.; Dalke, A.; Schulten, K. VMD: visual molecular dynamics. *Journal of Molecular Graphics* **1996**, 14(1), 33-38,  
 587 [https://doi.org/10.1016/0263-7855\(96\)00018-5](https://doi.org/10.1016/0263-7855(96)00018-5).  
 588  
 589  
 590

591 **Supplementary Table S1. Raw data for countries with cooccurring P681R, N501Y and D614G**  
 592 **mutations.** Frequencies for each region and continent are provided as sub-totals and totals respec-  
 593 tively.

Continent	Region	Country	Frequency	Sub-totals	Totals
Africa	Eastern Africa	Reunion	2	2	15
Africa	Sub-Saharan Africa	Botswana	2	13	
Africa	Sub-Saharan Africa	Malawi	2		
Africa	Sub-Saharan Africa	Mozambique	1		
Africa	Sub-Saharan Africa	Republic of the Congo	1		
Africa	Sub-Saharan Africa	South Africa	7		
Americas	Latin America and the Caribbean	Argentina	1	259	1078
Americas	Latin America and the Caribbean	Brazil	214		
Americas	Latin America and the Caribbean	Chile	2		
Americas	Latin America and the Caribbean	Colombia	1		
Americas	Latin America and the Caribbean	Czech Republic	23		
Americas	Latin America and the Caribbean	Ecuador	5		
Americas	Latin America and the Caribbean	Martinique	4		
Americas	Latin America and the Caribbean	Mexico	5		
Americas	Latin America and the Caribbean	Puerto Rico	4		
Americas	Northern America	Canada	5	819	
Americas	Northern America	USA	814		
Asia	Central Asia	Uzbekistan	1	1	856
Asia	Eastern Asia	Japan	4	5	
Asia	Eastern Asia	South Korea	1	11	
Asia	South-eastern Asia	Cambodia	2		
Asia	South-eastern Asia	Malaysia	1		
Asia	South-eastern Asia	Philippines	4		
Asia	South-eastern Asia	Singapore	3		
Asia	South-eastern Asia	Thailand	1		
Asia	Southern Asia	Bangladesh	5	28	
Asia	Southern Asia	India	21		
Asia	Southern Asia	Pakistan	1		
Asia	Southern Asia	Sri Lanka	1		
Asia	Western Asia	Georgia	3	811	
Asia	Western Asia	Iraq	1		
Asia	Western Asia	Israel	18		
Asia	Western Asia	Turkey	786		

Asia	Western Asia	United Arab Emirates	3		
Europe	Eastern Europe	Bulgaria	2	73	1728
Europe	Eastern Europe	Hungary	2		
Europe	Eastern Europe	Poland	12		
Europe	Eastern Europe	Romania	10		
Europe	Eastern Europe	Slovakia	45		
Europe	Eastern Europe	Ukraine	2		
Europe	Northern Europe	Denmark	69	472	
Europe	Northern Europe	Estonia	26		
Europe	Northern Europe	Finland	2		
Europe	Northern Europe	Iceland	8		
Europe	Northern Europe	Ireland	38		
Europe	Northern Europe	Latvia	1		
Europe	Northern Europe	Lithuania	2		
Europe	Northern Europe	Norway	2		
Europe	Northern Europe	Sweden	191		
Europe	Northern Europe	United Kingdom	133		
Europe	Southern Europe	Croatia	12	117	
Europe	Southern Europe	Greece	7		
Europe	Southern Europe	Italy	47		
Europe	Southern Europe	Malta	1		
Europe	Southern Europe	Portugal	28		
Europe	Southern Europe	Slovenia	1		
Europe	Southern Europe	Spain	21		
Europe	Western Europe	Austria	8	1066	
Europe	Western Europe	Belgium	21		
Europe	Western Europe	France	836		
Europe	Western Europe	Germany	133		
Europe	Western Europe	Luxembourg	8		
Europe	Western Europe	Netherlands	22		
Europe	Western Europe	Switzerland	38		
Oceania	Australia and New Zealand	New Zealand	1	1	1

594

595

596

597

598

599

600

601  
602  
603  
604  
605  
606  
607  
608  
609

**Supplementary Table S2. Observation of key mutations and their combinations on GISAID.** Live analysis from EpiCoV (as of 3 November 2021) containing the results for each of our MOCI, and counts for VOC/VOI according to GISAID. This is compared with data downloaded from GISAID (as of 1 November 2021) and analysed using our in-house alignment pipeline (complete human-origin genomes) after filtering for correct dates in the metadata ('YYYY-MM-DD').

	<b>GISAID live analysis using EpiCoV</b>	<b>Data downloaded from GISAID (correct dates)</b>
<b>D614G</b>	4,245,283	4,129,785
<b>N501Y</b>	1,255,861	1,220,340
<b>P681R</b>	1,887,719	1,844,450
<b>D614G + N501Y</b>	1,251,360	1,216,080
<b>D614G + P681R</b>	1,882,323	1,839,455
<b>N501Y + P681R</b>	3,728	3,688
<b>D614G + N501Y + P681R</b>	3,718	3,678
<b>D614G + N501Y + P681H</b>	1,112,530	1,082,482
<b>Alpha</b>	1,127,758	1,096,665
<b>Beta</b>	39,478	37,935
<b>Gamma</b>	96,217	92,164
<b>Delta</b>	1,871,122	1,828,769
<b>Alpha + P681R</b>	2,393	2,352 <sup>a</sup>
<b>Beta + P681R</b>	24	24 <sup>b</sup>
<b>Gamma + P681R</b>	294	294 <sup>c</sup>
<b>Delta + N501Y</b>	917	912 <sup>d</sup>
<sup>a</sup> mostly present in the sub-lineage Q.4 (2119 samples, since 13 December 2020), B.1.1.7 (230 samples, 17 October 2020) <sup>b</sup> mostly present in B.1.351 (23 samples, since 30 March 2021) <sup>c</sup> mostly present in sub-lineage P.1.8 (180 samples, since 20 August 2021) <sup>d</sup> mostly present in B.1.617.2 (476 samples, since 29 March 2021), sub-lineage AY.33 (155 samples, since 30 April 2021)		

610  
611  
612  
613  
614  
615

616

617

618

619

620

621

622

623

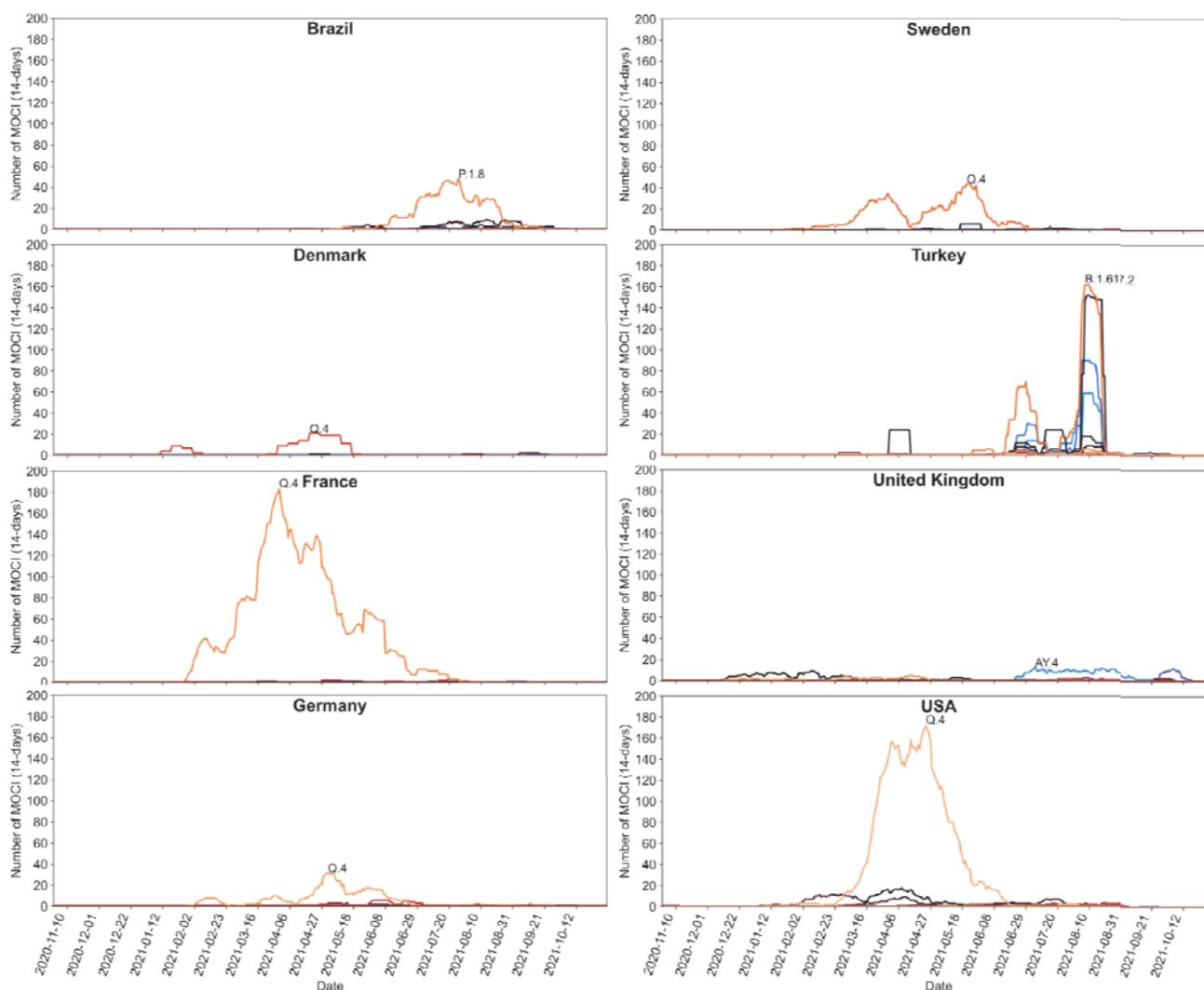
624

**Supplementary Table S3. Observation of other key SARS-CoV-2 Spike mutations.** Live analysis from EpiCoV (as of 3 November 2021) containing the results for other key Spike mutations according to GISAID. This is compared with data downloaded from GISAID (as of 1 November 2021) and analysed using our in-house alignment pipeline after filtering for correct dates in the metadata ('YYYY-MM-DD').

	GISAID live analysis using EpiCoV	Data downloaded from GISAID	Data downloaded from GISAID (correct dates)
<b>Total</b>	4,835,087	4,766,139	
<b>Human</b>	4,831,865	4,763,023	
<b>Human complete</b>	4,748,462	*4,294,469	4,177,098
<b>L18F</b>	212,488	197,955	191,238
<b>T20N</b>	108,445	93,507	89,541
<b>H69del</b>	1,123,439	1,138,302	1,108,661
<b>Y145del</b>	1,364	2,617	2,563
<b>A222V</b>	403,457	367,047	352,024
<b>K417N</b>	44,282	43,114	41,449
<b>N439K</b>	35,253	35,135	34,365
<b>L452R</b>	2,304,322	1,919,808	1,877,945
<b>Y453F</b>	1,255	1,242	1,221
<b>G476S</b>	933	799	764
<b>S477N</b>	70,642	69,479	67,649
<b>T478I</b>	902	892	879
<b>V483A</b>	296	264	243
<b>E484K</b>	225,663	205,807	197,301
<b>E484Q</b>	11,834	10,746	9,998
<b>E780Q</b>	5,931	5,733	5,627
<b>V1176F</b>	125,546	108,949	103,509
*complete human-origin genomes			

625

626



**Supplementary Figure S4. Frequency of the P681R, N501Y and D614G mutations cooccurring in the eight countries with more than 50 observations between October 2020 and October 2021. For visual clarity, only the SARS-CoV-2 lineage with the maximum count is labelled.**

640  
641  
642  
643  
644

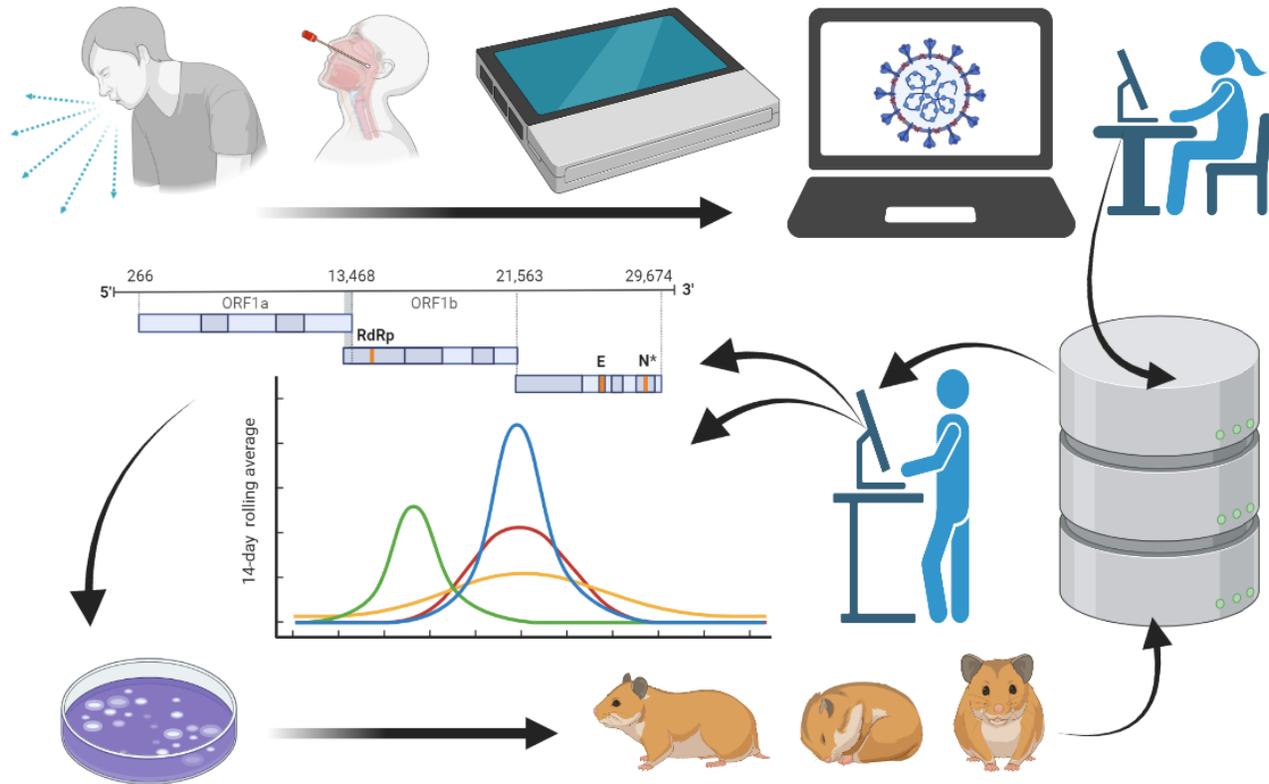
**Supplementary Table S5. GISAID data on key SARS-CoV-2 mutations classified by VOC/VOI.** Live analysis from EpiCoV (as of 3 November 2021) containing the results for each of our MOCI, and counts for VOC/VOI according to GISAID. This is compared with data downloaded from GISAID (as of 1 November 2021) and analysed using our in-house alignment pipeline after filtering for correct dates in the metadata ('YYYY-MM-DD').

	<b>GISAID live analysis using EpiCoV</b>	<b>Data downloaded from GISAID</b>	<b>Data downloaded from GISAID (correct dates)</b>
<b>Total</b>	4,835,087	4,766,139	
<b>Human origin</b>	4,831,865	4,763,023	
<b>Human and complete</b>	4,748,462	*4,294,469	4,177,098
<b>D614G</b>	4,692,072	4,245,283	4,129,785
<b>N501Y</b>	1,282,756	1,255,861	1,220,340
<b>P681R</b>	2,288,985	1,887,719	1,844,450
<b>Alpha</b>	1,135,705	1,127,758	1,096,665
<b>Beta</b>	38,865	39,478	37,935
<b>Gamma</b>	113,091	96,217	92,164
<b>Delta</b>	2,267,806	1,871,122	1,828,769
<b>Lambda</b>	8,771	8,149	8,084
<b>Mu</b>	13,042	12,447	11,814
*complete human-origin genomes			

645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660

661

### Graphical Abstract

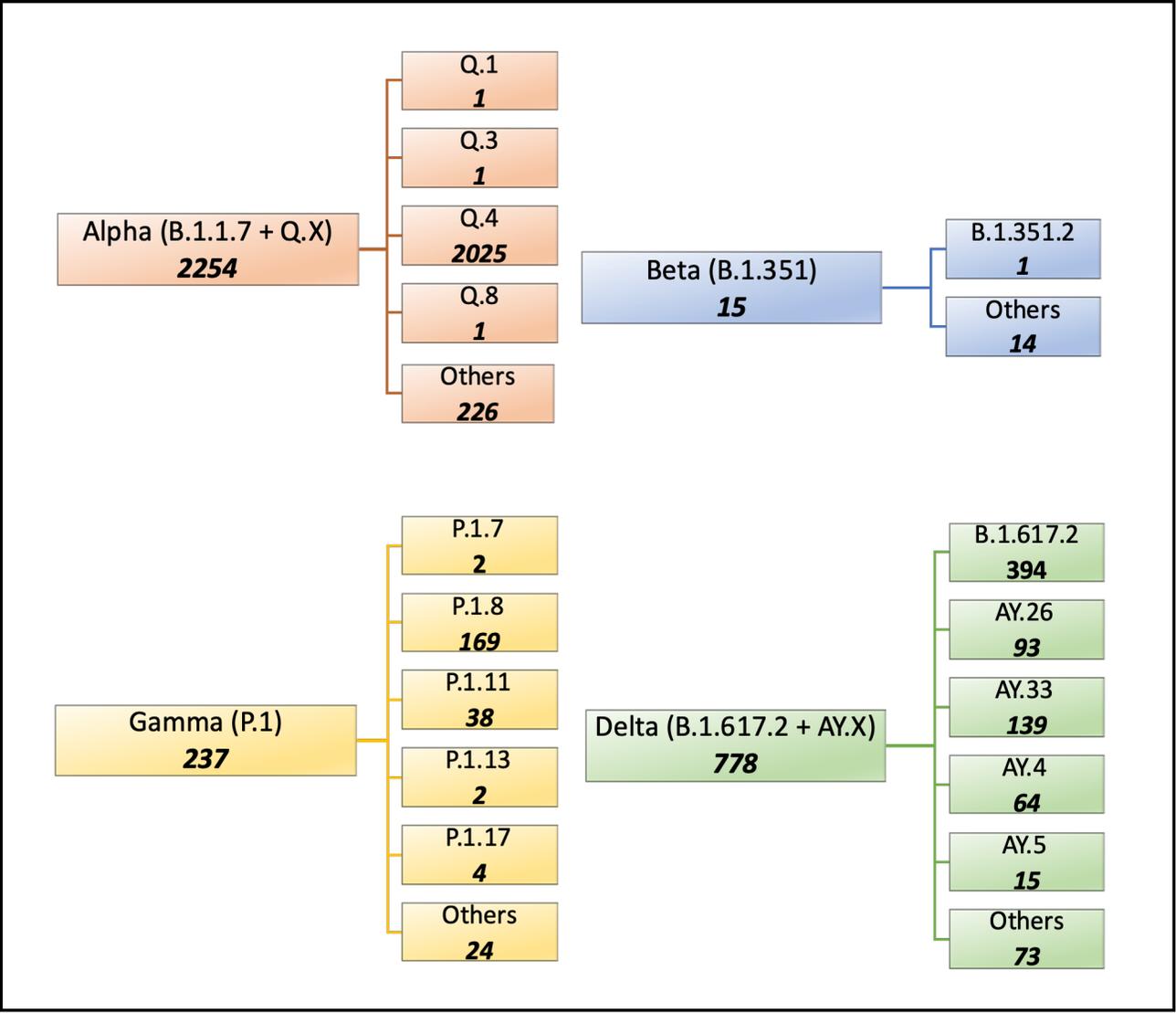


## Cooccurrence of N501Y, P681R and other key mutations in SARS-CoV-2 Spike

Created in [BioRender.com](https://www.biorender.com)

662

663



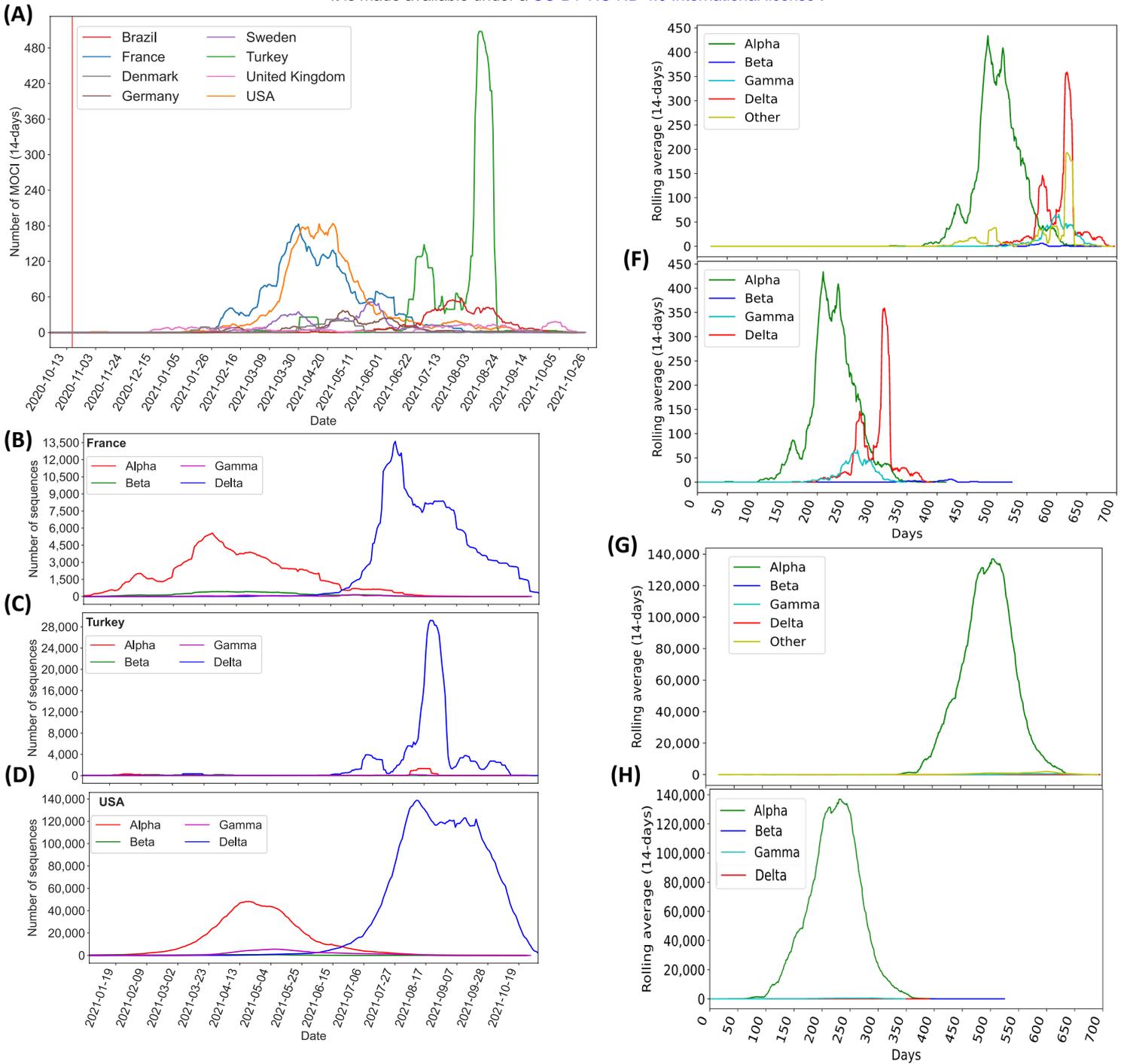
VOI  
Mu (B.1.621)  
**12**

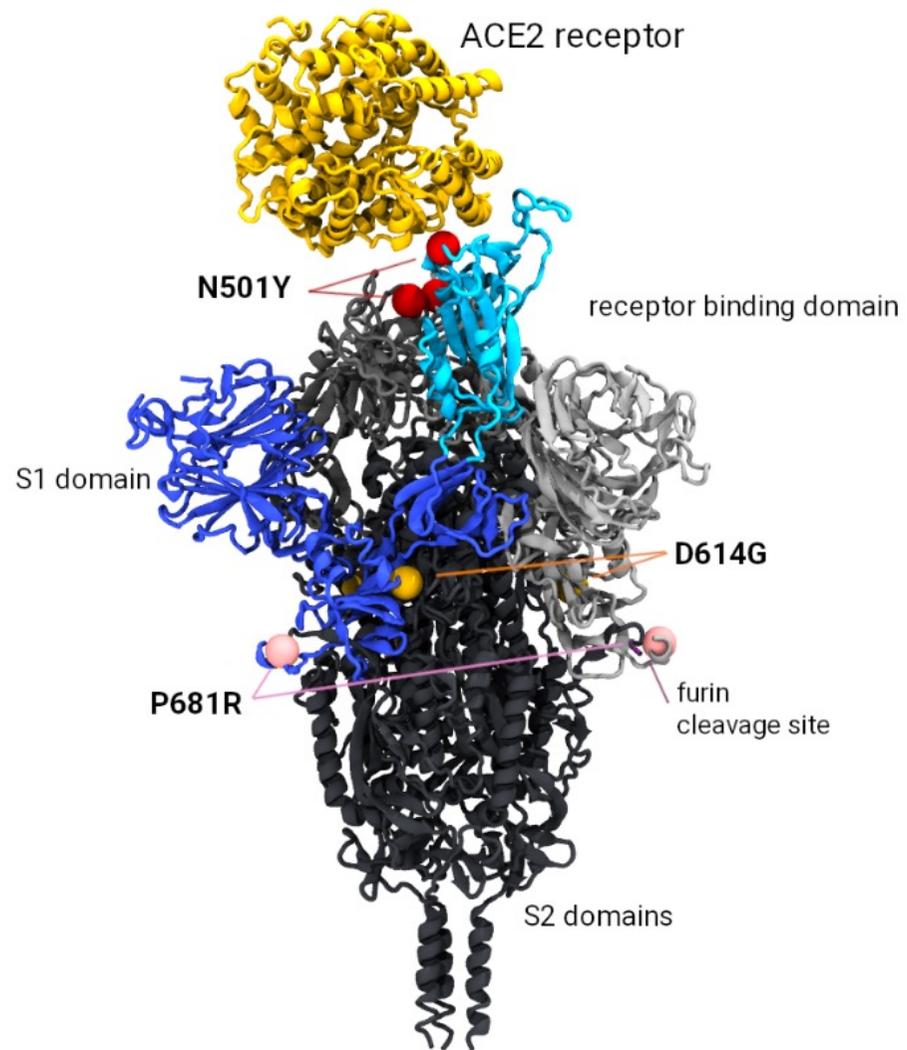
Former VOI  
Kappa (B.1.617.1)  
**8**

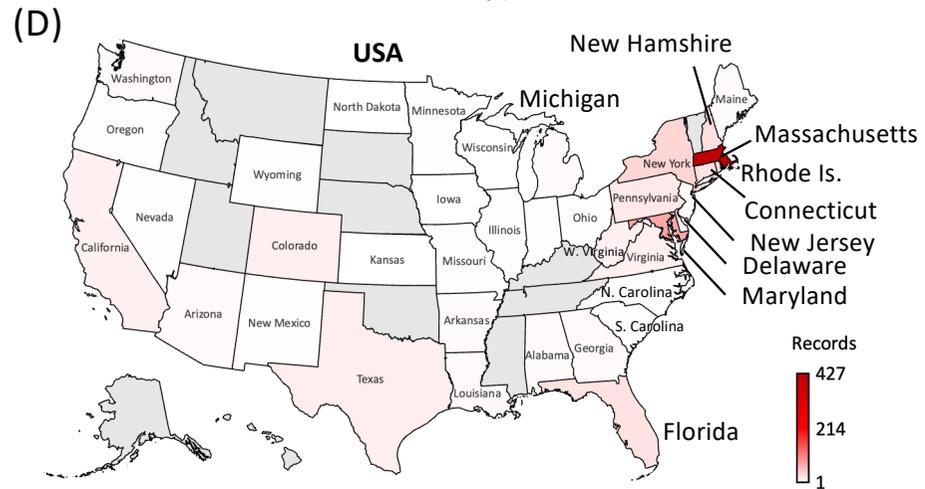
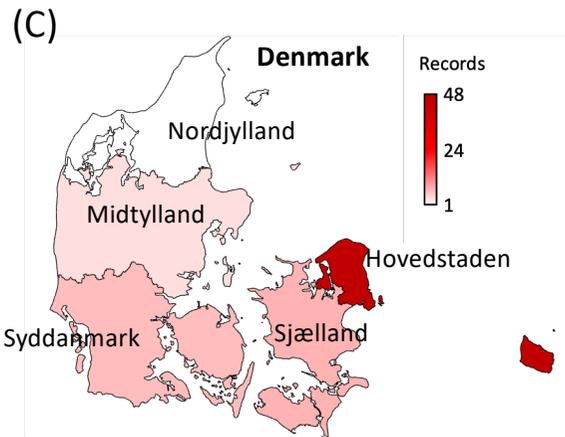
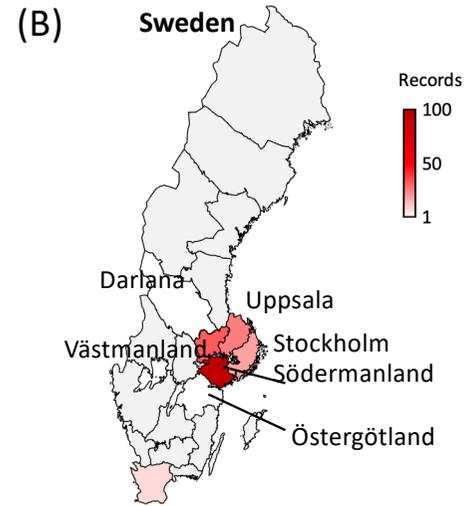
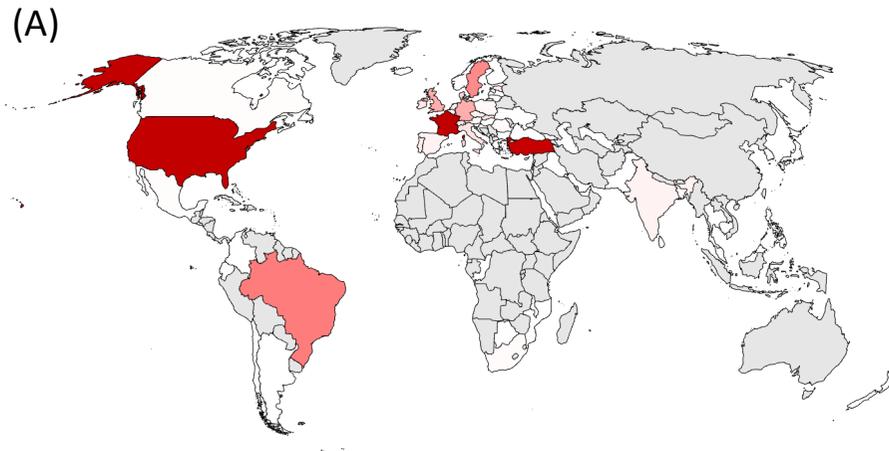
Former VOI  
Iota (B.1.596)  
**1**

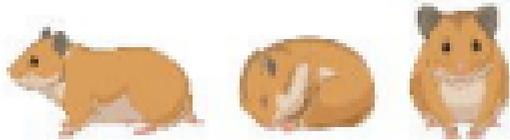
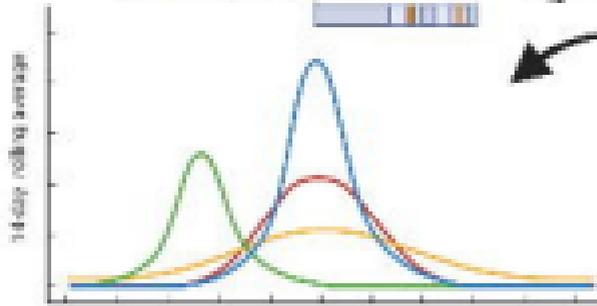
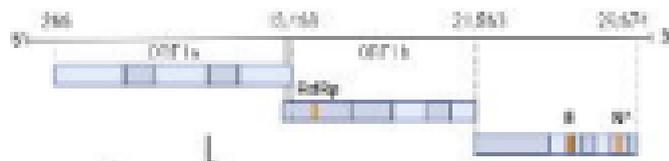
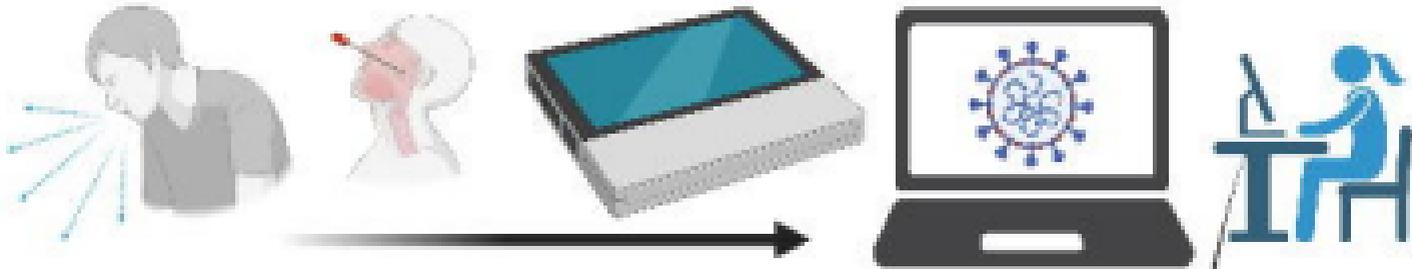
Others  
**373**











Cocurrence of N501Y, P681R and other key mutations in SARS-CoV-2 Spike