

1 **Title: Computation of Antigenicity Predicts SARS-CoV-2 Vaccine Breakthrough**
2 **Variants**

3 **Authors:**

4 Ye-fan Hu^{1,3,†}, Jing-chu Hu^{2,†}, Hua-rui Gong^{1,†}, Antoine Danchin^{1,5}, Ren Sun¹, Hin
5 Chu⁴, Ivan Fan-Ngai Hung³, Kwok Yung Yuen⁴, Kelvin Kai-Wang To⁴, Bao-zhong
6 Zhang^{2,*}, Thomas Yau^{3,*}, Jian-Dong Huang^{1,2,*}

7 **Affiliations:**

8 ¹ School of Biomedical Sciences, Li Ka Shing Faculty of Medicine, University of
9 Hong Kong, 3/F, Laboratory Block, 21 Sassoon Road, Hong Kong, China

10 ² CAS Key Laboratory of Quantitative Engineering Biology, Shenzhen Institute of
11 Synthetic Biology, Shenzhen Institutes of Advanced Technology, Chinese Academy
12 of Sciences, Shenzhen 518055, China

13 ³ Department of Medicine, Li Ka Shing Faculty of Medicine, University of Hong
14 Kong, 4/F Professional Block, Queen Mary Hospital, 102 Pokfulam Road, Hong
15 Kong, China

16 ⁴ Department of Microbiology, Li Ka Shing Faculty of Medicine, University of Hong
17 Kong, 19/F T Block, Queen Mary Hospital, 102 Pokfulam Road, Hong Kong, China

18 ⁵ Kodikos Labs / Stellate Therapeutics, Institut Cochin, 24 rue du Faubourg Saint-
19 Jacques, 75014 Paris, France

20

21 [†] These authors contributed equally to this work.

22 * Corresponding authors. B.Z.Z. (bz.zhang3@siat.ac.cn), T.Y. (tyaucc@hku.hk),
23 J.D.H. (jdhuang@hku.hk)

24

25 **Abstract**

26 It has been reported that multiple SARS-CoV-2 variants of concerns (VOCs)
27 including B.1.1.7 (Alpha), B.1.351 (Beta), P.1 (Gamma), and B.1.617.2 (Delta) can
28 reduce neutralisation by antibodies, resulting in vaccine breakthrough infections.
29 Virus-antiserum neutralisation assays are typically performed to monitor potential
30 vaccine breakthrough strains. However, such experimental-based methods are slow
31 and cannot instantly validate whether newly emerging variants can break through
32 current vaccines or therapeutic antibodies. To address this, we sought to establish a
33 computational model to predict the antigenicity of SARS-CoV-2 variants by sequence
34 alone and in real time. In this study, we firstly identified the relationship between the
35 antigenic difference transformed from the amino acid sequence and the antigenic
36 distance from the neutralisation titres. Based on this correlation, we obtained a
37 computational model for the receptor binding domain (RBD) of the spike protein to
38 predict the fold decrease in virus-antiserum neutralisation titres with high accuracy
39 (~0.79). Our predicted results were comparable with experimental neutralisation titres
40 of variants, including B.1.1.7 (Alpha), B.1.351 (Beta), B.1.617.2 (Delta), B.1.429
41 (Epsilon), P.1 (Gamma), B.1.526 (Iota), B.1.617.1 (Kappa), and C.37 (Lambda), as
42 well as SARS-CoV. Here, we firstly predicted the fold of decrease of B.1.1.529
43 (Omicron) as 17.4-fold less susceptible to neutralisation. We visualised all 1521
44 SARS-CoV-2 lineages to indicate variants including B.1.621 (Mu), B.1.630, B.1.633,
45 B.1.649, and C.1.2, which can induce vaccine breakthrough infections in addition to
46 reported VOCs B.1.351 (Beta), P.1 (Gamma), B.1.617.2 (Delta), and B.1.1.529
47 (Omicron). Our study offers a quick approach to predict the antigenicity of SARS-
48 CoV-2 variants as soon as they emerge. Furthermore, this approach can facilitate
49 future vaccine updates to cover all major variants. An online version can be accessed
50 at <http://jdlab.online> .

51

52 Up to January 2022, there have been several SARS-CoV-2 variants including B.1.1.7
53 (Alpha)¹⁻⁵, B.1.351 (Beta)^{2,3,6,7}, P.1 (Gamma)^{1,2,8}, and B.1.617.2 (Delta)^{9,10} that are
54 experimentally tested to lead vaccine breakthrough infections, thus they have been
55 designated as variants of concerns (VOCs) by the world health organization (WHO).
56 There is a concern that other untested emerging variants may lead to vaccine
57 breakthrough infections¹¹⁻¹⁶. The most recent case is the validation of B.1.1.529
58 (Omicron). The current virological and epidemiological techniques took several
59 weeks to validate whether the variant is capable of reducing the efficacy of current
60 vaccines^{17,18} or therapeutic antibodies^{18,19}, even though their viral sequences have
61 been shared in real time via the Global Initiative for Sharing All Influenza Data
62 (GISAID)²⁰. The speed of validation of vaccine breakthrough variants can hardly
63 catch up with the fast-emerging rate of new variants. Thus, it is crucial to develop new
64 approaches for identifying the next potential vaccine breakthrough variant as soon as
65 it is reported.

66 Here, we established a computational approach for predicting the antigenicity of
67 SARS-CoV-2 variants from viral sequences alone, with the aim to accelerate the
68 identification of potential vaccine breakthrough variants. Our approach is founded on
69 the concept of antigenic mapping, also named antigenic cartography. This method has
70 been used to monitor vaccine breakthrough variants of influenza virus using
71 haemagglutination inhibition (HI) assay data^{21,22}, dengue virus²³ and SARS-CoV-2
72 circulating strains²⁴ using pairwise antisera data. In antigenic mapping, the antigenic
73 distance is calculated from the fold change of the neutralisation titre between the
74 reference virus and its variant, to measure the change of antigenicity between two
75 variants. A computational approach for predicting antigenic distances to indicate
76 vaccine breakthrough variants could theoretically provide much more rapid results
77 once the variant sequence is reported. Past studies proposed a linear relationship
78 between amino acid changes in antigenic sites and neutralisation fold decrease²⁵⁻²⁹.
79 Computational prediction approaches based on such a relationship could also provide
80 reliable estimates of neutralisation titres for existing antiserum against the vaccine
81 breakthrough variants with similar accuracy to experiment-based approaches used in
82 previous studies²⁵⁻²⁹. However, these predictions were optimised for influenza virus
83 instead of SARS-CoV-2. For example, the neutralisation titre decrease of any SARS-
84 CoV-2 variant should be less than that of SARS-CoV comparing to the ancestral

85 strain of SARS-CoV-2, because the cross protection between the SARS-CoV-2
86 variant and the ancestral strain is stronger than that between SARS-CoV and SARS-
87 CoV-2. Thus, it is difficult to use a linear relationship to predict the decrease in
88 neutralisation titre which saturates with the increase in the mutation numbers of
89 variants. A SARS-CoV-2 optimised model for predicting antigenicity is urgently
90 needed.

91 In this study, we established a computational sequence-based method to predict the
92 antigenicity of SARS-CoV-2 variants to reveal potential vaccine breakthrough
93 variants. This method can also predict the neutralisation titre of VOCs in comparison
94 to the ancestral strain of SARS-CoV-2. Our predicted results were comparable with
95 experimental neutralisation titres of VOCs, including B.1.1.7 (Alpha), B.1.351 (Beta),
96 B.1.617.2 (Delta), B.1.429 (Epsilon), P.1 (Gamma), B.1.526 (Iota), B.1.617.1 (Kappa),
97 and C.37 (Lambda), as well as SARS-CoV. Here, we predicted that B.1.1.529
98 (Omicron) is 17.4-fold less susceptible to neutralisation, which is consistent with
99 reported decrease folds ranging from 10 to 40^{17,18}.

100 **A computational model for predicting antigenicity of SARS-CoV-2 variants**

101 To predict the antigenicity of SARS-CoV-2 variants, we firstly integrated the reported
102 conformational or linear epitopes (**Fig. S1 & Table S1**) on the SARS-CoV-2 Spike
103 protein (**Fig. 1a**) with the reported experimental virus-antiserum neutralisation titres
104 against SARS-CoV-2 variants including B.1.1.7¹⁻⁵, B.1.351^{2,3,6,7}, and P.1^{1,2,8} (**Table**
105 **S2a**). Considering the distinct assays used in the different studies, we standardised the
106 neutralisation titres of each variant to the titre of the ancestral strain of SARS-CoV-2
107 (lineage A) using the same assay in each study on a log₂ scale, and thus we got
108 observed antigenic distance (H_{ab}) from neutralisation titres (**Fig. 1b**). For the antigenic
109 difference (D_{ab}), we used Poisson distance to represent the difference between two
110 amino acid sequences (**Fig. 1b**). By comparing the observed antigenic distance with
111 the antigenic difference, we found a relationship between observed antigenic distance
112 and the antigenic difference: $H_{ab}=T_{max}\cdot D_{ab}/(D_{50}+D_{ab})$, where T_{max} is the maximal fold
113 of decrease and D_{50} is the antigenic difference which may lead to neutralisation
114 decrease at the 50% level of the maximal decrease (the fold change between SARS-
115 CoV-2 and SARS-CoV). This relationship described that the decrease of
116 neutralisation titre increases with the accumulation of amino acid changes, and then

117 reaches at the maximal decrease (**Figs. 1c-d**). Based on this correlation, we obtained a
118 computational model using the receptor binding domain (RBD) of the spike protein to
119 predict the fold decrease in virus-antiserum neutralisation titres with higher accuracy
120 (**~0.79, the calculation of accuracy in Methods**) compared with other fragments of
121 spike (entire spike, N terminal domain plus RBD, or S1, **Fig. 1d**). With repeated 5-
122 fold or 10-fold cross validation (**Fig. 1d**), we found that prediction using RBD is
123 relatively robust in terms of root-mean-square error (RMSE), mean absolute error
124 (MAE), coefficient of determination (R^2) and accuracy.

125 To further validate our model, we predicted the fold decreases in neutralisation titres
126 (comparing to the ancestral of SARS-CoV-2) of multiple variants including B.1.1.7
127 (Alpha), B.1.351 (Beta), B.1.617.2 (Delta), B.1.429 (Epsilon), P.1 (Gamma), B.1.526
128 (Iota), B.1.617.1 (Kappa), and C.37 (Lambda), as well as SARS-CoV and WIV1-CoV
129 using datasets without the variant that we aimed to validate. Previous studies have
130 reported that VOCs can elicit vaccine breakthrough infections, which correlated with
131 fold decreases in the neutralisation titres from experimental assays was disclosed
132 (**Table S2**). Our predicted results were highly consistent with the neutralisation assay
133 results (**Fig. 1e**). We also predicted the fold of decrease in neutralisation titre of the
134 most recent VOC, B.1.1.529 (Omicron). Considering 15 mutations in the spike of
135 B.1.1.529 (Omicron), the variant is estimated to have a 17.44-fold (95% confidence
136 interval: 13.7, 22.2) decrease in neutralisation titre (shown as a blue point in **Fig. 1c**). ,
137 The predicted result is consistent with reported decrease folds ranging from 10 to 40
138 ^{17,18}. This result alarmed the risk of vaccine breakthrough or re-infection of B.1.1.529
139 (Omicron) due to the dramatic decrease in neutralization.

140 **The prediction of potential vaccine breakthrough strains**

141 To predict the next potential SARS-CoV-2 vaccine breakthrough variants, we
142 visualised the antigenicity of all available SARS-CoV-2 variants as an indicator of
143 their vaccine breakthrough potential. We firstly selected all 1521 lineage variants
144 using PANGO ³⁰ updated on December 6, 2021 (**Table S3**) to predict their
145 antigenicity. Then we calculated the pairwise distances of different variants. For
146 visualising these results, we captured two principal components from the high-
147 dimensional data of antigenic distance ²⁵. We used all spike amino acid sequences to
148 plot the ‘genetic map’ of SARS-CoV-2 to represent the genetic difference among

149 different variants (**Fig. 2a-b**). We then plotted the ‘antigenic map’ using the predicted
150 antigenic distances (**Fig. 2c-d**, online versions available at <http://jdlab.online>).

151 Based on the relationship between neutralisation titre fold change and protective
152 efficacy³¹, it was convenient to set up some ‘cut-offs’ in the current vaccine coverage.
153 We included phase 3 and real-world results of vaccine efficacy or effectiveness, as
154 well as neutralisation titre data from phase 1 and 2 studies (**Table S4-5**). Thus, we got
155 the relationship between neutralisation titre and protective efficacy against a
156 symptomatic COVID-19 (**Fig. 2e**). A 3.93-fold decrease in neutralisation titres
157 induced by VOCs that can dampened the efficacy of some vaccines to lower than 50%.
158 In this way, one cut-off of 1.98 arbitrary units (A.U.) represented a 3.93-fold decrease
159 in the neutralisation titre (shown as a pink circle in **Fig. 2c-d**). All variants outside this
160 cut-off have the potential to be vaccine breakthrough variants. By comparing the
161 “genetic map” and antigenic map, we can set up the border of antigenic map.
162 Although there are >200 mutations in the SARS-CoV and WIV1-CoV spike (**Fig 2a**),
163 the antigenic distance is around 4.9 A.U. which mean ~ 30-fold decrease in the
164 neutralisation titre (shown as a dark red circle in **Fig. 2c-d**).

165 To reveal the distribution of variant, we plotted the density of variants on the ‘genetic
166 map’ and antigenic map due to overlapping dots. In the genetic map, hotspots are
167 located at lineage A (>10%) and B.1 (>40%) mainly, as well as AY.* and P.1 (**Fig.**
168 **2b**). While in the antigenic map, hotspots are placed at lineage A (>40%) mainly,
169 together with AY.* (**Fig. 2d**). Although most variants were shown to be close to the
170 ancestral strain (**Figs. 2b&d**), multiple variants were found to decrease neutralisation
171 titres significantly (**Fig. 2c**). In addition to reported VOCs including B.1.351 (Beta,
172 containing sub-lineages like B.1.351.2 and B.1.351.5)^{2,3,6,7}, P.1 (Gamma, containing
173 sub-lineages like P.1.11 and P.1.3)^{1,2,8}, B.1.617.2 (Delta, containing sub-lineages
174 AY.*)⁹, and B.1.621 (Mu, containing sub-lineage B.1.621.1), B.1.1.529 (Omicron)
175 showed over 3.93-fold decrease in the neutralisation titre. Other variants B.1.630,
176 B.1.633, B.1.649, and C.1.2 also have the potential to be vaccine breakthrough
177 variants with more than 3.93-fold decrease (**Fig. 2c**). Besides the pandemic of
178 B.1.617.2 (Delta)⁹ and the outbreak of B.1.1.529 (Omicron), multiple variants should
179 be investigated immediately as they have the potential to become tomorrow’s VOCs.

180 **Discussion**

181 Predicting neutralisation responses against all SARS-CoV-2 variants based on
182 sequences alone is vital for selecting the next vaccine seeds for the development of
183 effective COVID-19 vaccines. We established a computational approach to predict
184 neutralisation titres and validated these predictions using experimental data. Our
185 computational approach could potentially provide the first hints of whether a newly
186 identified variant can break through vaccines just by its sequence information, which
187 would greatly shorten the time for the crucial early warning of emerging vaccine
188 breakthrough strains.

189 In the prediction of the antigenicity of SARS-CoV-2 variants, we proposed that the
190 limit of neutralisation titre decrease is set by SARS-CoV (**Fig. 1**). In recent studies,
191 SARS-CoV is ~ 36-fold less susceptible to neutralisation comparing to the ancestral
192 strain of SARS-CoV-2. Based on this result, a non-linear curve was established to
193 describe the relationship between the observed antigenic distance and the antigenic
194 difference. We further performed calculation using different fragments of the Spike
195 protein (**Fig. 1d**). Among the Spike protein and the RBD, NTD-RBD, and S1
196 fragments, we found the prediction using amino acid sequences of RBD was able to
197 estimate the neutralisation titre more accurately than the others (**Figs. 1d**). Thus, we
198 used the RBD-based computations to determine the neutralisation titres.

199 A major concern of our computation of the neutralisation titre is that the data is based
200 on diverse neutralisation assays of serum samples from both patients and vaccinees
201 against both live virus and pseudovirus (**Table S2**). Although the results were
202 consistent qualitatively, the variation of fold change is too large to be ignored (**Fig.**
203 **1e**). Considering the variation in the real world, we set up values 2-fold or less than
204 the experimental values as the criteria based on previous studies²⁸. It is better to
205 establish a convenient and standardised neutralisation pipeline in the future, like the
206 haemagglutination inhibition (HI) assay for influenza virus. Such a pipeline can allow
207 the precise estimation of neutralisation titres. Together with estimating the association
208 of neutralisation with protection, it will help to develop next generation vaccines.

209 It is crucial to update vaccines to cover all vaccine breakthrough strains that have
210 significant amino acid and glycosylation changes to prevent further infectious
211 outbreaks. However, not all predicted SARS-CoV-2 vaccine breakthrough variants
212 will have the chance to cause an outbreak due to their changed viral fitness³² or by

213 pure luck. Based on previous studies of influenza viruses, it is possible for variants to
214 have alterations that change the antigenicity, but fail to cause outbreaks in the wider
215 population ³³. Considering immune escape elicited by variants, updating current
216 vaccine seeds with new variants should extend the vaccine coverage. As SARS-CoV-
217 2 showed different variant directions in the antigenic map (**Fig. 2**), the use of multiple
218 virus seeds based on the different directions might be appropriate to cover all major
219 variants in the long term. Our method could help in the selection of SARS-CoV-2
220 variants for updating vaccines.

221 **Methods**

222 **Antigenic footprint**

223 We collected 149 confirmed conformational epitopes with protein structures released
224 in the Protein Data Bank (PDB) (<https://www.rcsb.org/>) or annotated epitope
225 footprints and 76 linear epitopes published in the literature (**Table S1**). We plotted the
226 footprint of all Spike protein epitopes from the aforementioned 225 epitopes using R-
227 3.6.6.

228 **Antigenic distances from neutralisation data**

229 We calculated antigenic distances from the neutralisation data based on previous
230 publications²⁶. For virus variant a , reference virus b , and antiserum β (referencing
231 virus b), we defined the antigenic distance of variant a to reference virus b in terms of
232 the standardised log titre as $H_{ab} = \log_2 T_{a\beta} - \log_2 T_{b\beta}$, where $T_{b\beta}$ is the titre of antiserum β
233 against virus b , and $T_{a\beta}$ is the titre of antiserum β against virus a ²⁶. Merged data with
234 reference virus lineage A (the ancestral strain of SARS-CoV-2) were collected from
235 several publications (**Table S2**).

236 **Genetic and antigenic difference calculation**

237 We selected 1521 SARS-CoV-2 lineages using PANGO (v.3.1.15) updated on
238 December 6, 2021 (<https://cov-lineages.org/>). Spike protein amino acid sequences of
239 these lineages were obtained from GISAID, using the earliest collected for each
240 lineage (**Table S3**). All sequences with neutralisation titres were also included (**Table**
241 **S3**). For genetic distances, we used Molecular Evolutionary Genetics Analysis
242 (MEGA) X to calculate the pairwise distances among Spike protein amino acid
243 sequences in the SARS-CoV-2 variants using a Poisson model. For antigenic distance,
244 we used an information theory-based approach *p-all-epitope*^{27,28} to measure the
245 pairwise distances among amino acid sequences of the antigenic footprint ('antigenic
246 positions'). The distance is based on the number of different amino acids n_d between
247 two n -mer viral sequences of variants a and b . Under the assumption that the number
248 of amino acid substitutions per site follows a Poisson distribution, we can then
249 calculate the distance between a and b as $D_{ab} = -\ln(1 - n_d/n)$.

250 **Modelling and performance measurement**

251 A model considering the maximal neutralisation titre decrease was applied to examine
252 the antigenic distance from the neutralisation data H_{ab} and our computed results D_{ab} as
253 $H_{ab}=T_{max}\cdot D_{ab}/(D_{50}+D_{ab})$, where T_{max} is the maximal decrease and D_{50} is the antigenic
254 difference which may lead to neutralisation decrease at the 50% level of the maximal
255 decrease. The predicted neutralisation titre is then given as
256 $P_{ab}\approx\hat{H}_{ab}=T_{max}\cdot D_{ab}/(D_{50}+D_{ab})$. Root-mean-square error (RMSE), mean absolute error
257 (MAE), and coefficient of determination (R-squared R^2) were used to measure the
258 performance of the linear correlation.

259 Reproducibility was determined by pairwise sequences and neutralisation titres.
260 Neutralisation titre data were converted into variables by calculating the relative
261 difference in the neutralisation titres between reference virus and variant against the
262 antiserum. Accuracy was the percentage of correctly predicted neutralisation titres
263 using amino acid sequences. Based on previous studies²⁸, computational values 2-
264 fold or less than the experimental values were considered to be similar (correct) and
265 those more than 2-fold lower were considered dissimilar (error). Here, 10-time
266 repeated 5-fold and 10-fold cross validation were applied in terms of root-mean-
267 square error (RMSE), mean absolute error (MAE), coefficient of determination (R-
268 squared R^2), and accuracy.

269 **Genetic and antigenic maps**

270 After calculating genetic and antigenic distances, we used classical multidimensional
271 scaling (CMDs) to display the data as a plot using R-3.6.6. We set up SARS-CoV-2
272 lineage A as the origin and scaled the data in two and three dimensions. We then
273 acquired the genetic and antigenic maps of SARS-CoV-2 lineages. An online version
274 can be obtained at <http://jdlab.online> .

275 **Logistic model**

276 Following past studies³¹, we used a logistic model in R-3.6.6 to describe the
277 relationship between antigenic distance (neutralization level) and protective
278 efficacy/effectiveness: $E=1/(1+\exp(-k(H-H_{50})))$. E is the protective
279 efficacy/effectiveness at a specific neutralization level H . H is the mean of
280 neutralisation titres in vaccinees divided by corresponding mean of titres in
281 convalescent patients, which is the antigenic distance to convalescent patients in log 2.

282 H_{50} is the antigenic distance at which an individual will have a 50% protective
283 efficacy/effectiveness.

284 **References**

- 285 1. Garcia-Beltran WF, Lam EC, Denis KS, et al. Multiple SARS-CoV-2 variants escape
286 neutralization by vaccine-induced humoral immunity. *Cell* 2021.
- 287 2. Chen RE, Zhang X, Case JB, et al. Resistance of SARS-CoV-2 variants to
288 neutralization by monoclonal and serum-derived polyclonal antibodies. *Nature Medicine* 2021.
- 289 3. Wang P, Nair MS, Liu L, et al. Antibody Resistance of SARS-CoV-2 Variants
290 B.1.351 and B.1.1.7. *Nature* 2021.
- 291 4. Xie X, Liu Y, Liu J, et al. Neutralization of SARS-CoV-2 spike 69/70 deletion,
292 E484K and N501Y variants by BNT162b2 vaccine-elicited sera. *Nature Medicine* 2021.
- 293 5. Muik A, Wallisch A-K, Sanger B, et al. Neutralization of SARS-CoV-2 lineage
294 B.1.1.7 pseudovirus by BNT162b2 vaccine-elicited human sera. *Science* 2021: eabg6105.
- 295 6. Supasa P, Zhou D, Dejnirattisai W, et al. Reduced neutralization of SARS-CoV-2
296 B.1.1.7 variant by convalescent and vaccine sera. *Cell* 2021.
- 297 7. Zhou D, Dejnirattisai W, Supasa P, et al. Evidence of escape of SARS-CoV-2 variant
298 B.1.351 from natural and vaccine induced sera. *Cell* 2021.
- 299 8. Liu Y, Liu J, Xia H, et al. Neutralizing Activity of BNT162b2-Elicited Serum. *New*
300 *England Journal of Medicine* 2021.
- 301 9. Edara V-V, Lai L, Sahoo MK, et al. Infection and vaccine-induced neutralizing
302 antibody responses to the SARS-CoV-2 B.1.617.1 variant. *bioRxiv* 2021: 2021.05.09.443299.
- 303 10. Wall EC, Wu M, Harvey R, et al. Neutralising antibody activity against SARS-CoV-2
304 VOCs B.1.617.2 and B.1.351 by BNT162b2 vaccination. *The Lancet* 2021; **397**(10292):
305 2331-3.
- 306 11. Hacısuleyman E, Hale C, Saito Y, et al. Vaccine Breakthrough Infections with SARS-
307 CoV-2 Variants. *New England Journal of Medicine* 2021.
- 308 12. COVID-19 breakthrough case investigations and reporting. 2021.
- 309 13. Kustin T, Harel N, Finkel U, et al. Evidence for increased breakthrough rates of
310 SARS-CoV-2 variants of concern in BNT162b2 mRNA vaccinated individuals. *medRxiv*
311 2021: 2021.04.06.21254882.
- 312 14. Jacobson KB, Pinsky BA, Rath MEM, et al. Post-vaccination SARS-CoV-2
313 infections and incidence of the B.1.427/B.1.429 variant among healthcare personnel at a
314 northern California academic medical center. *medRxiv* 2021: 2021.04.14.21255431.
- 315 15. Abu-Raddad LJ, Chemaitelly H, Butt AA. Effectiveness of the BNT162b2 Covid-19
316 Vaccine against the B.1.1.7 and B.1.351 Variants. *New England Journal of Medicine* 2021.
- 317 16. Cavanaugh AM, Fortier S, Lewis P, et al. COVID-19 Outbreak Associated with a
318 SARS-CoV-2 R.1 Lineage Variant in a Skilled Nursing Facility After Vaccination Program -
319 Kentucky, March 2021. *MMWR Morb Mortal Wkly Rep* 2021; **70**(17): 639-43.
- 320 17. Carreño JM, Alshammary H, Tcheou J, et al. Activity of convalescent and vaccine
321 serum against SARS-CoV-2 Omicron. *Nature* 2021.
- 322 18. Liu L, Iketani S, Guo Y, et al. Striking Antibody Evasion Manifested by the Omicron
323 Variant of SARS-CoV-2. *Nature* 2021.
- 324 19. Cao Y, Wang J, Jian F, et al. Omicron escapes the majority of existing SARS-CoV-2
325 neutralizing antibodies. *Nature* 2021.
- 326 20. Elbe S, Buckland-Merrett G. Data, disease and diplomacy: GISAID's innovative
327 contribution to global health. *Global Challenges* 2017; **1**(1): 33-46.
- 328 21. Smith DJ, Lapedes AS, de Jong JC, et al. Mapping the Antigenic and Genetic
329 Evolution of Influenza Virus. *Science* 2004; **305**(5682): 371-6.
- 330 22. Koel BF, Burke DF, Bestebroer TM, et al. Substitutions Near the Receptor Binding
331 Site Determine Major Antigenic Change During Influenza Virus Evolution. *Science* 2013;
332 **342**(6161): 976-9.
- 333 23. Katzelnick LC, Fonville JM, Gromowski GD, et al. Dengue viruses cluster
334 antigenically but not as discrete serotypes. *Science* 2015; **349**(6254): 1338-43.
- 335 24. Liu C, Ginn HM, Dejnirattisai W, et al. Reduced neutralization of SARS-CoV-2
336 B.1.617 by vaccine and convalescent serum. *Cell* 2021; **184**(16): 4220-36.e13.
- 337 25. Sitaras I. Antigenic Cartography: Overview and Current Developments. In:

- 338 Spackman E, ed. *Animal Influenza Virus: Methods and Protocols*. New York, NY: Springer
339 US; 2020: 61-8.
- 340 26. Neher RA, Bedford T, Daniels RS, Russell CA, Shraiman BI. Prediction, dynamics,
341 and visualization of antigenic phenotypes of seasonal influenza viruses. *Proceedings of the*
342 *National Academy of Sciences* 2016; **113**(12): E1701-E9.
- 343 27. Pan K, Subieta KC, Deem MW. A novel sequence-based antigenic distance measure
344 for H1N1, with application to vaccine effectiveness and the selection of vaccine strains.
345 *Protein Engineering, Design and Selection* 2011; **24**(3): 291-9.
- 346 28. Anderson CS, McCall PR, Stern HA, Yang H, Topham DJ. Antigenic cartography of
347 H1N1 influenza viruses using sequence-based antigenic distance calculation. *BMC*
348 *Bioinformatics* 2018; **19**(1): 51.
- 349 29. Gupta V, Earl DJ, Deem MW. Quantifying influenza vaccine efficacy and antigenic
350 distance. *Vaccine* 2006; **24**(18): 3881-8.
- 351 30. O'Toole Á, Scher E, Underwood A, et al. Assignment of epidemiological lineages in
352 an emerging pandemic using the pangolin tool. *Virus Evolution* 2021; **7**(2).
- 353 31. Khoury DS, Cromer D, Reynaldi A, et al. Neutralizing antibody levels are highly
354 predictive of immune protection from symptomatic SARS-CoV-2 infection. *Nature Medicine*
355 2021.
- 356 32. Hie B, Zhong ED, Berger B, Bryson B. Learning the language of viral evolution and
357 escape. *Science* 2021; **371**(6526): 284-8.
- 358 33. Plante JA, Liu Y, Liu J, et al. Spike mutation D614G alters SARS-CoV-2 fitness.
359 *Nature* 2020.
- 360 34. Woo H, Park S-J, Choi YK, et al. Developing a Fully Glycosylated Full-Length
361 SARS-CoV-2 Spike Protein Model in a Viral Membrane. *The Journal of Physical Chemistry*
362 *B* 2020; **124**(33): 7128-37.

363

364 **Acknowledgements**

365 We acknowledge the authors and originating and submitting laboratories of the
366 sequences from GISAID's EpiCoV Database on which this research is based
367 (GISAID acknowledgments are in Table S3).

368 **Funding:** The work was supported by grants from the Health and Medical Research
369 Fund, the Food and Health Bureau, The Government of the Hong Kong Special
370 Administrative Region (COVID190117, COVID1903010) and Guangdong Science
371 and Technology Department (2020B1212030004) to J.H. J.H. thanks the L & T
372 Charitable Foundation and the Program for Guangdong Introducing Innovative and
373 Entrepreneurial Teams (2019BT02Y198) for their support.

374 **Author contributions:** Y.F.H., B.Z.Z., T.Y., H.C., K.K.W.T., and J.D.H. designed
375 the study. Y.F.H. and J.C.H. analysed the sequences from GISAID and performed the
376 computations and built the online tool. I.F.N.H. and K.K.W.T. recruited the patients
377 and volunteers and collected samples from patients and volunteers. performed the
378 neutralisation assay. Y.F.H., A.D., R.S., K.Y.Y., K.T., H.R.G, and J.D.H. analysed

379 the results. Y.F.H., A.D., T.Y., and J.D.H. wrote the initial draft, and all authors
380 edited the final version.

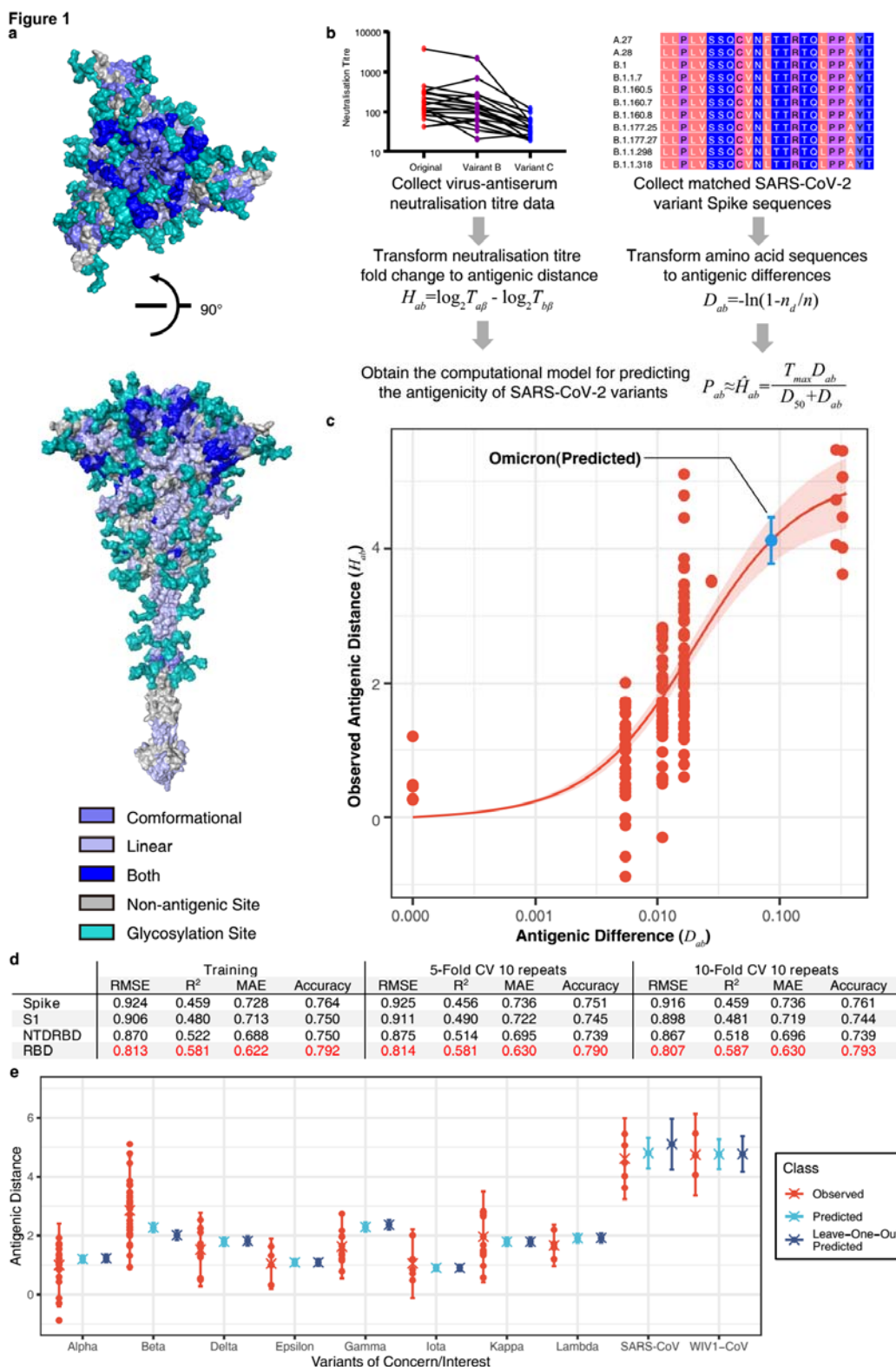
381 **Competing interests:** All authors declare no competing interests.

382 **Data and materials availability:** All sequence data listed in TableS3 are from
383 GISAID's EpiCoV Database.

384 **Supplementary Materials**

385 **Figures S1 to S2**

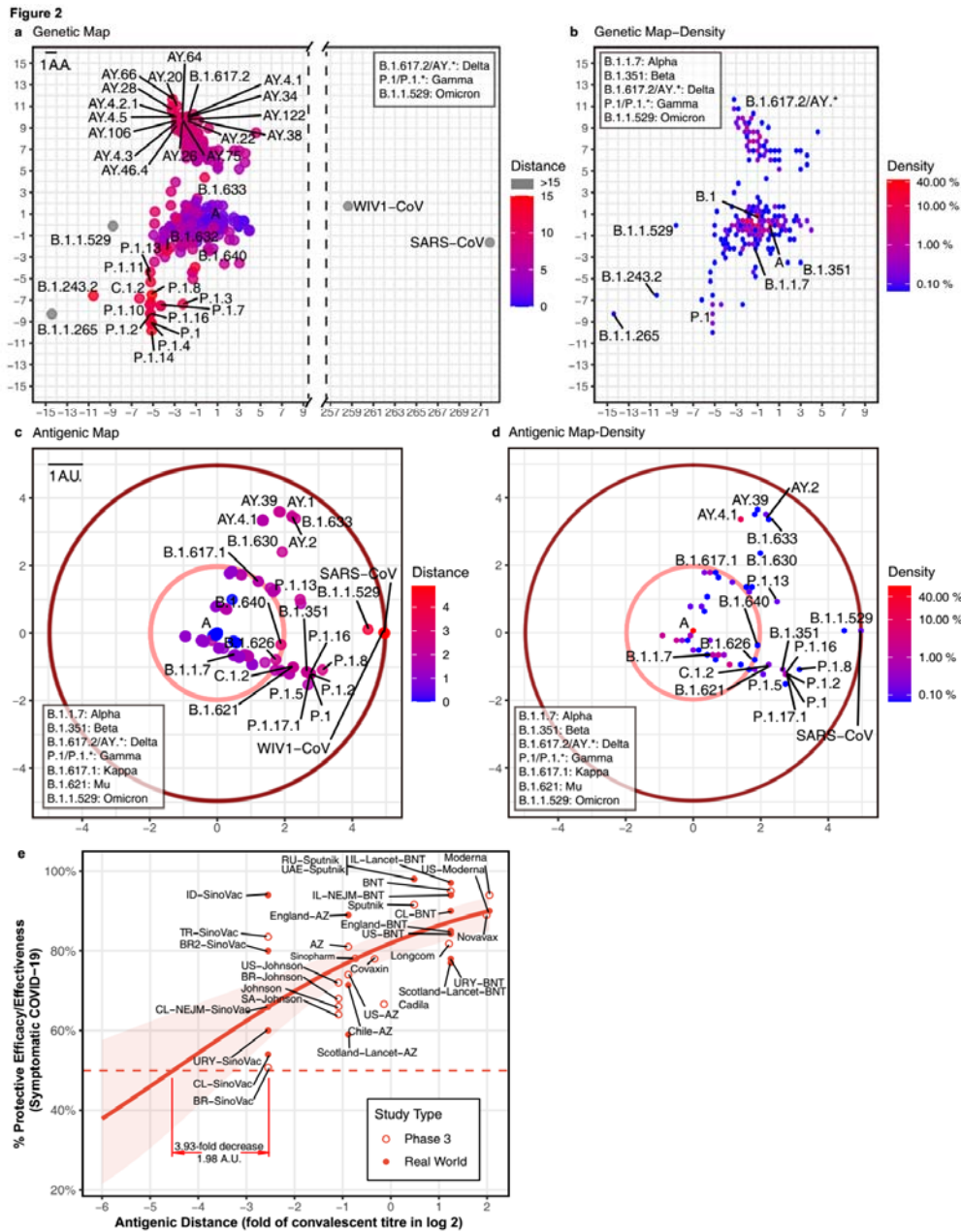
386 **Tables S1 to S5**



387

388 **Fig. 1| Sequence-based prediction of antigenic distance.** (a) The top view and the
 389 side view of antigenic sites on the full-length Spike protein³⁴. The conformational
 390 epitopes are coloured in slate and linear epitopes in light blue. Some antigenic

391 positions in both conformational epitopes and linear epitopes are coloured in blue. All
392 glycosylation sites are in teal. **(b)** A flowchart of the process to establish the sequence-
393 based computational model of SARS-CoV-2 antigenicity. The antigenic distance of
394 variant a to reference virus b from neutralisation titre was defined as $H_{ab} = \log_2 T_{a\beta} -$
395 $\log_2 T_{b\beta}$, where β , $T_{a\beta}$, and $T_{b\beta}$ denote antiserum (referencing virus b), the titre of
396 antiserum β against virus b , and the titre of antiserum β against virus a ²⁶. The
397 antigenic distance of variant a to reference virus b from amino acid sequences was
398 defined as $D_{ab} = -\ln(1 - n_d/n)$, where n_d is the number of amino acid substitutions
399 between variant a and reference virus b , n is the number of antigenic sites. Then, we
400 proposed a relationship between observed antigenic distance and the antigenic
401 difference: $H_{ab} = T_{max} \cdot D_{ab} / (D_{50} + D_{ab})$, where T_{max} is the maximal fold of decrease and
402 D_{50} is the antigenic difference which may lead to neutralisation decrease at the 50%
403 level of the maximal decrease. **(c)** The relationship between the antigenic difference
404 and the observed antigenic distance. The predicted antigenic distance of B.1.1.529
405 (Omicron) is marked in cyan. **(d)** The performance of the model in different
406 fragments of the spike protein in terms of root-mean-square error (RMSE), mean
407 absolute error (MAE), coefficient of determination (R-squared R^2), and accuracy. **(e)**
408 Predicted versus observed antigenic distances of variants of concern. Here, The
409 observed antigenic distances as fold decreases in the neutralisation titres of variants of
410 concern versus the original strain on a log 2 scale. Each point shows the mean of
411 antigenic distances in each assay. Predicted antigenic distances are based on the
412 prediction in (c). Leave-one-out predicted antigenic distances are predicted based on
413 the datasets without the variant that we aim to compare.



414

415 **Fig. 2| Genetic and antigenic mapping of SARS-CoV-2 variants.** (a) Genetic map
 416 of SARS-CoV-2 variant strains shows amino acid mutation numbers of spike proteins,
 417 and (b) the density of genetic map shows distribution of variants. The vertical and
 418 horizontal axes represent the measured relative genetic distances (1 amino acid/1 A.A.
 419 = 1 amino acid difference). (c) Antigenic map of SARS-CoV-2 variant strains shows
 420 the antigenic distance between variants, and (d) the density of antigenic map shows
 421 distribution of variants. Variants outside the pink circle are vaccine breakthrough
 422 candidates. The red circle suggested the border of antigenic map. The antigenic

423 distance is based on RBD amino acid sequences. The vertical and horizontal axes
424 represent the measured relative antigenic distances (1 arbitrary unit/1 A.U. = 1-fold
425 decrease in the neutralisation titre on a log 2 scale). Colours show the antigenic
426 distance to the SARS-CoV-2 original strain (lineage A). (e) Relationship between
427 antigenic distance (mean of neutralisation titres in vaccinees divided by corresponding
428 mean of titres in convalescent patients in log 2) and protection from SARS-CoV-2
429 infection. The reported mean neutralization level from phase 1 or 2 studies (**Table S4**)
430 and the protective efficacy or effectiveness from phase 3 trials or real-world studies
431 (**Table S5**) for different vaccines. The red line indicates the logistic model, and the
432 red shading indicates the 95% confidence interval of the model. Here, we mark the
433 basis of setting up the cut-off of 3.93-fold decrease (1.98 A.U.).

434

435

Figure 1

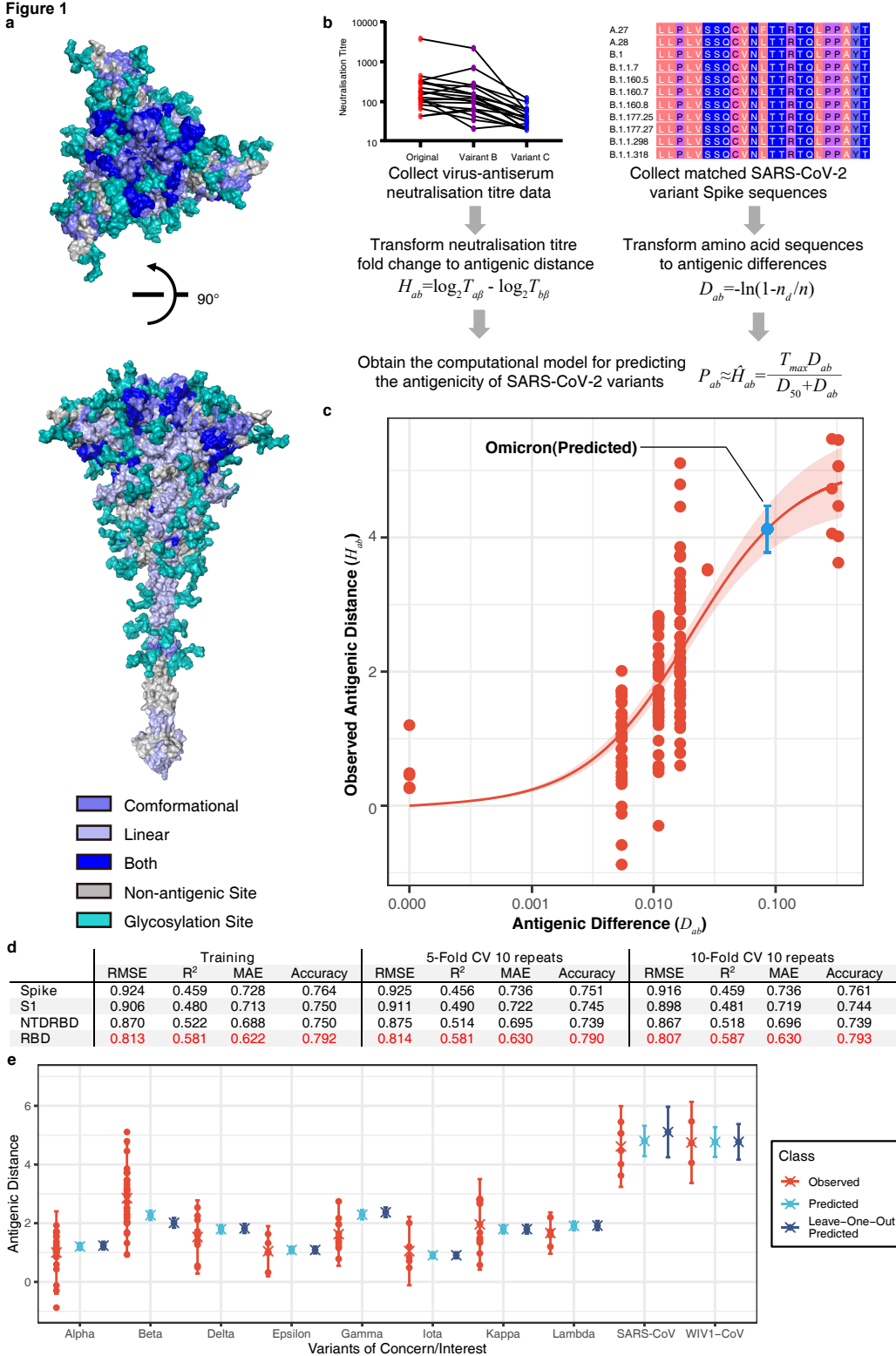
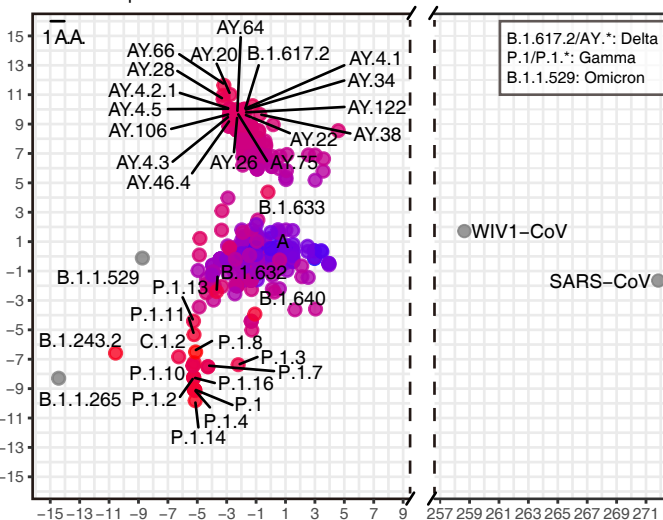
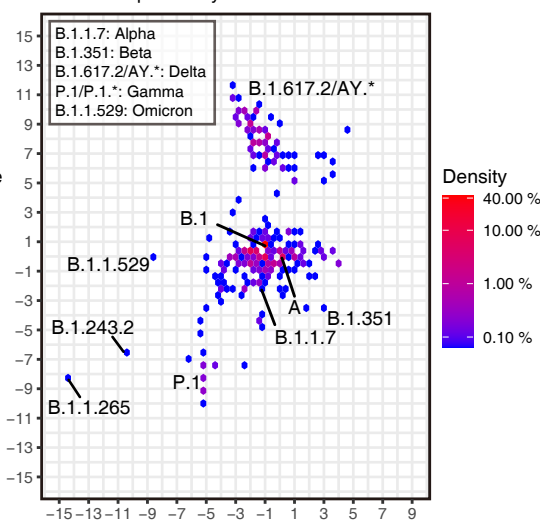
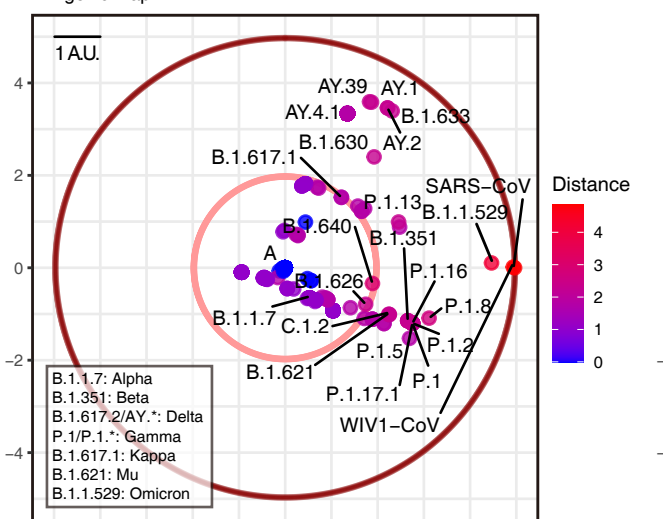


Figure 2**a Genetic Map****b Genetic Map-Density****c Antigenic Map****d Antigenic Map-Density**