

Predicting the antigenic evolution of SARS-COV-2 with deep learning

Wenkai Han^{1,2,#}, Ningning Chen^{1,2,#}, Xinzhou Xu^{3,4,#}, Adil Sahil^{1,2}, Juexiao Zhou^{1,2}, Zhongxiao Li^{1,2}, Huawen Zhong², Ruochi Zhang⁵, Yu Wang⁵, Shiwei Sun^{6,7,*}, Peter Pak-Hang Cheung^{3,4,*},
Xin Gao^{1,2,*}

¹Computer Science Program, Computer, Electrical and Mathematical Sciences and Engineering Division, King Abdullah University of Science and Technology (KAUST), Thuwal 23955-6900, Kingdom of Saudi Arabia

²Computational Bioscience Research Center, King Abdullah University of Science and Technology, Thuwal 23955-6900, Kingdom of Saudi Arabia

³Department of Chemical Pathology, Faculty of Medicine, Chinese University of Hong Kong, Hong Kong, China

⁴Li Ka Shing Institute of Health Sciences, Chinese University of Hong Kong, Hong Kong, China

⁵Syneron Technology, Guangzhou, 510000, China

⁶Key Lab of Intelligent Information Processing, Institute of Computing Technology, Chinese Academy of Sciences, Beijing, 100190, China

⁷University of Chinese Academy of Sciences, Beijing 100049, China

[#]The first three authors contributed equally to this paper.

*All correspondence should be addressed to X.G. (xin.gao@kaust.edu.sa), P.C. (ppcheung@cuhk.edu.hk), and S.S. (dwsun@ict.ac.cn).

24 **Abstract**

25 The severe acute respiratory syndrome coronavirus 2 (SARS-COV-2) antigenic profile evolves in response to
26 the vaccine and natural infection-derived immune pressure, resulting in immune escape and threatening public
27 health. Exploring the possible antigenic evolutionary potentials improves public health preparedness, but it is
28 limited by the lack of experimental assays as the sequence space is exponentially large. Here we introduce the
29 Machine Learning-guided Antigenic Evolution Prediction (MLAEP), which combines structure modeling, multi-
30 task learning, and genetic algorithm to model the viral fitness landscape and explore the antigenic evolution via
31 *in silico* directed evolution. As demonstrated by existing SARS-COV-2 variants, MLAEP can infer the order of
32 variants along antigenic evolutionary trajectories, which is also strongly correlated with their sampling time. The
33 novel mutations predicted by MLAEP are also found in immunocompromised covid patients and newly emerging
34 variants, like XBB1.5. The predictions of MLAEP were validated by conducting in vitro neutralizing antibody
35 binding assay, which demonstrated that the model-generated variants displayed significantly increased immune
36 evasion ability compared with the controls. In sum, our approach enables profiling existing variants and
37 forecasting prospective antigenic variants, thus may help guide the development of vaccines and increase
38 preparedness against future variants. Our model is available at <https://mlaep.cbrc.kaust.edu.sa>.

39

40 **Introduction**

41 As the number of infection cases increased and the virus spread globally, novel mutations in the virus genome
42 emerged¹⁻⁴. At the time of April 2022, there are more than one million variants in the virus genome identified
43 and uploaded to the Global Initiative on Sharing Avian Influenza Database (GISAID). The mutations often
44 implicate the changes to the SARS-COV-2 properties³. Although most mutations decrease the virulence and
45 transmissibility of the virus⁵, some individual or combinatorial mutations substantially improve the
46 transmissibility with enhanced cell entry efficacy⁶, or ablate the neutralizing antibodies response elicited by
47 infection or vaccine^{1,7}, resulting in high-risk variants. For example, the Alpha (B.1.1.7) variant of concern (VOC)
48 spread worldwide through a higher human ACE2 binding affinity and transmissibility than the original Wuhan
49 strain⁸. The Beta and the Gamma lineage abolished the neutralizing antibodies elicited by approved COVID-19
50 vaccines⁹. The Delta variant became a dominant strain worldwide with the increased transmissibility and
51 mortality^{10,11}. Recently, the heavily mutated Omicron variant caused new waves due to the extremely high rate
52 of spread and the ability to evade the double-vaccinated person¹².

53 A substantial fraction of neutralizing antibodies, including monoclonal antibodies and those induced by the
54 vaccines, target the spike receptor-binding domain (RBD)¹³⁻¹⁵. Antibodies targeting the RBD have been divided
55 into four categories according to their binding epitopes¹⁶. Class 1 and class 2 antibodies bind the surface of the
56 receptor-binding motif (RBM) and thus compete with ACE2 for RBD binding. Mutations in the RBM region, in
57 turn, decreased neutralization by these antibodies. Class 3 antibodies bind the opposite side of the receptor-
58 binding motif, contain less overlap with the ACE2-binding footprint, provide the potential for synergistic effects
59 when combined with Class 1 and 2 antibodies for intercepting ACE2 binding¹⁷. Class 4 antibodies target a highly
60 conserved region among sarbecoviruses and thus are generally more resistant to the variants¹⁸. However, the
61 emerging viral lineages such as Omicron and BA.2 can still lead to a substantial loss of neutralization¹⁹.

62 Understanding the role of the mutations and how they are linked to transmissibility and immune escape are thus
63 of great importance. There have been an expanding set of analyses characterizing these problems^{5,18,20-23}. Starr
64 et al.⁵ and Greaney et al.¹⁸ performed deep mutational scanning (DMS) on the entire Spike RBD sequences of
65 SARS-COV-2 on the yeast surface to determine the impact of single-position substitutions on the binding ability
66 to ACE2 and monoclonal antibodies. These assayed experiments provide a unique resource for understanding

67 the properties of variants. However, the wet-lab experiments are resource and time-consuming, and cannot be
68 scaled to the large protein sequence space. Maher et al.²³ characterized the potential risks of the single-position
69 substitutions with a computational model and forecasted the driver mutations that may appear in emerging VOCs.
70 Despite their effectiveness in modeling the risks at the single-mutant level, the newly emerging VOCs (e.g.,
71 Delta, Omicron) often possess multiple mutations in the RBD region, which directly influences the ACE2
72 binding and antibody escape. For example, the Omicron variant contains 15 mutations in the RBD region and
73 obtains considerable antigenic escape ability²⁴. Moreover, the effects of mutations are context-dependent, such
74 that the epistatic interactions among the mutations limit the application scenario of the single-mutant-based
75 methods²⁵.

76 The sequence space of protein variants grows exponentially when multiple mutations are considered, while
77 measuring the functionality of the variant sequences far exceeding the capacity of wet-lab experiments. Machine
78 learning methods have been proposed for solving the problem²⁶⁻²⁸. Alexander et al.²⁹ trained a large-scale
79 transformer model with the self-supervised protein language modeling objective, while the model can infer the
80 effects of mutations without supervision. Chloe et al.³⁰ combined linear regression with the Potts model, resulting
81 in a data-efficient variant fitness inference model. These models have been proven to be effective in the protein
82 engineering field for inferring the fitness landscape of proteins.

83 Inspired by these tools, Hie et al.²⁰ showed that language models trained on a set of evolutionarily related
84 sequences are capable of predicting the potential risks of COVID variants with multiple mutations, and Karim
85 et al.²² further combined the language model score with structural modeling to monitor the risks of existing
86 variants. These computational tools can work as high-risk variant monitors and help us predict the risks of the
87 emerging variants. However, as these methods focus on prediction and rely on existing data, they do not provide
88 detailed views for ‘perspective’ variants and antigenic evolutionary potential. Taft et al.²¹ performed deep
89 learning on the RBM sequences and built a predictive profile for the COVID variants in ACE2 binding and
90 antibody escape for class 1, 2, and 3 antibodies. The proposed framework works quite well in finding prospective
91 mutations, but they still have limitations: the mutations are found by brute-force search, so they only focused on
92 a small subset of the RBD region, missed a large part of the Class 3 antibody epitopes and did not take the class
93 4 antibodies into consideration.

94 In this work, we presented the MLAEP, built upon the existing data and approaches to forecast the combinatorial
95 mutations in the entire RBD region that contains high antigenic evolutionary potential and may occur in the
96 future. We hypothesized that under high immune pressure, the virus would tend to escape the antibody
97 neutralization over a short-term time scale, and therefore the forecasting problem transforms into a search
98 problem: starting from an initial sequence, it searches for a variant sequence within some edit distance range that
99 has an improved antibody escape potential without losing much ACE2 binding ability. With the DMS datasets
100 that directly measure the binding affinity of RBD variants towards ACE2 and eight antibodies from four classes,
101 we built a multi-task deep learning model that could simultaneously predict the binding/escaping specificity of
102 the variants towards the ACE2 and eight antibodies. Furthermore, we used existing variants with their sampling
103 date from the GISAID database to validate our hypothesis: we found a surprisingly high correlation between our
104 model scores and the variants' sampling time (Spearman $r=0.65$, $p<1e-308$). Next, with our model as the scoring
105 function, we used the genetic algorithm^{31, 32} to generate synthetic RBD variants with high ACE2 binding and
106 antibody escape potential. Interestingly, the *in silico* directed evolution shares similar mutations with the adaptive
107 evolution in immunocompromised COVID-19 patients³³⁻³⁵ and newly emerging variants like XBB.1.5. Finally,
108 we conducted *in vitro* neutralizing antibody binding assay to verify the ability of MLAEP to accurately forecast
109 variants with high immune evasion potential.

110

111 **Results**

112 **Overview of MLAEP**

113 We first developed and trained a multi-task deep neural network model capable of predicting the variant RBD
114 binding specificity towards the ACE2 and antibodies from four classes, as shown in Figure 1. The model receives
115 two inputs: the variant RBD sequences and the ACE2/antibody 3D structures, and outputs the binding
116 specificities of the two inputs. The model is then trained with a multi-task objective function to predict the
117 binding specificities of the variant sequences towards all targets simultaneously.

118 We fine-tuned the ESM-1b (evolutionary scale modeling) language model²⁹ for the sequence feature extraction.
119 The model is pre-trained on ~27 million nature protein sequences in the UniRef50 database³⁶. Fine-tuning the
120 model has been proven to be effective for a broad range of downstream tasks, including biophysical properties

121 prediction, structure prediction, and mutation effects prediction. With the ESM-1b model, the amino acid
122 sequences are converted into a dense vector representation. For the ACE2/antibodies structures, we first
123 transformed the 3D structures into graphs based on their contact maps and biophysical properties, then used the
124 structured transformer³⁷ for the structural feature extraction. With the two models as feature extraction modules,
125 we added nine parallel linear classification layers to learn the sequence to function mapping conditioned on the
126 binding target structures (Fig. 1a). As we have multiple binding targets for the variants, we used a hard-parameter
127 sharing scheme to perform multi-task learning, where all modules share the same parameters across all nine
128 tasks. Then, we trained the entire framework in an end-to-end manner. Finally, the model learns how to predict
129 binding specificity for ACE2 and eight antibodies. Given an input RBD variant sequence, our model outputs
130 nine scores corresponding to the ACE2 and eight antibodies. We defined the average of eight antibody scores as
131 the predicted antibody escaping potential.

132 Our key hypothesis is based on the antigenic evolution: the future viral variants tend to have a higher antibody
133 escaping potential without losing much ACE2 binding ability under high immune pressure. Thus, the
134 antibody/ACE2 binding specificity learned by our model can be used to provide a meaningful direction in
135 searching for novel variants that may cause future concern. Inspired by the progress in the machine learning-
136 guided protein engineering field^{26,27}, we used the trained multi-task model as the scoring function (Fig. 1a), took
137 the average prediction scores from all nine tasks as the fitness score, and used a modified genetic algorithm for
138 searching for novel variants with improved fitness (Fig. 1b). The genetic algorithm is inspired by the process of
139 natural selection, which iteratively evolves a group of candidates towards better fitness. The population of each
140 iteration is called a generation. In each generation, the fitness of the candidate sequences is evaluated with the
141 trained model. Then we filtered the populations by selecting the ones with higher fitness with higher probabilities
142 for breeding the next generation (Fig. 1c). Random mutations and crossover are also introduced to better explore
143 the search space. Genetic algorithm is known for performing well in solving combinatorial optimization
144 problems, thus fitting our needs in searching for novel variants. More details can be found in Methods.

145

146

147 **The effectiveness of the multi-task learning model**

148 MLAEP follows the machine learning-guided directed evolution paradigm, while the quality of generated
149 sequences largely depends on the sequence to function model. First, we validated the generalization ability of
150 the models to newly seen variants with 5-fold cross-validation. We collected and cleaned nine deep mutational
151 datasets containing 19132 variant sequences and their corresponding binding specificities towards ACE2 and
152 eight antibodies from four functional classes. (Methods) We then compared a range of models specifically
153 designed for protein engineering and assessed their classification performance in classifying the binders and
154 non-binders (Methods, Extended Data Fig. 1) from the variant sequences, including the augmented Potts³⁰ model,
155 the global UniRep³⁸ model, the eUniRep²⁶ model, the convolutional neural network (CNN), the long short
156 memory neural network (LSTM), the recurrent neural network (RNN), the linear regression model, the support
157 vector machine (SVM), the random forest and our model. We used 5-fold cross-validation to evaluate the
158 performance of all models. The dataset is imbalanced regarding the number of positive and negative samples for
159 all nine tasks. Thus, we reported the macro precision, macro recall, and macro F1 score to add more weights to
160 the minor classes. Combined with the structure features, our model outperforms the other advanced methods in
161 predicting the effects of mutations in all nine tasks (Fig. 2a, Supplementary Fig. 1, 2, Supp Table 1). As a result,
162 we focused on our model in the downstream analysis. We also performed ablation study for our model to show
163 the importance of each module. We found that both the fine-tuning step and the structure representations
164 improves the overall model performance (Extended Data Fig. 2). We also conducted external validation
165 experiments using several deep mutational scanning datasets^{39,40} in addition to variant RBDs, and found that our
166 model performed comparably and consistently well across all tasks (Supp Table 2).

167 To further validate the model's predictions for immune escape, we used the *in vitro* pseudovirus neutralization
168 test (pVNT) datasets⁴¹ that measured the cross-neutralizing effect⁴¹ of 17 RBD monoclonal antibodies against
169 pseudoviruses expressing the Spike protein of selected variants of concern (VOCs). The pVNT assay reported
170 the observed fold change in the IC₅₀ of the antibody response for these VOC-derived pseudoviruses, with lower
171 fold change score indicating greater immune evasion compared to the wild type (Wuhu-1) reference pseudovirus.
172 Across all pseudoviruses and antibodies tested, we found surprisingly high positive correlations (Fig. 2b,
173 Supplementary Fig. 3, Supp Table 3) between the predicted antibody escape potential and the log fold change in
174 the IC₅₀.

175 The Evo-velocity⁴² enables the inference of evolutionary dynamics for proteins with a deep learning model. It
176 was built upon the premise that global evolution occurs through local amino acid changes and leveraged protein
177 language models to model the local rules of evolution (Methods). We next assessed our model's ability in
178 inferring the evolutionary trajectory of the existing RBD sequences using the Evo-velocity. We used the existing
179 SARS-COV-2 RBD sequences from the GISAID database across a timescale of around 27 months, from Dec.
180 2019 to Mar. 2022. The existing GISAID variant sequences were first transformed into embeddings with our
181 multi-task model. On top of the embeddings, we assigned directions among them based on the changes in the
182 average score predicted by our model, which forms the evolutionary "vector field". We visualized the
183 embeddings in the two-dimensional space with the Uniform Manifold Approximation and Projection (UMAP)⁴³
184 (Methods). The variants of concern, including Alpha, Beta, Delta, and Omicron, were mapped into different
185 clusters, and the velocities among these variants matched well with the known evolutionary trajectory (Fig. 3a).
186 Despite the model being trained only with the RBD sequences, the pseudo time inferred with our model had a
187 Spearman correlation of 0.55 ($p < 1e-308$) with the known variant sampling time (Extended Data Fig. 3a,
188 Supplementary Fig. 4). While using the ESM-1b (the Evo-velocity default setting) model, the score dropped to
189 -0.38 ($p = 1.05e-243$) quickly (Extended Data Fig. 3b, c). We noted that a large set of mutations occur outside the
190 RBD region; this may explain the weak correlation between the ESM-1b model pseudo time and the sampling
191 time. Longer sequence length (e.g., using the entire Spike protein region) would lead to better performance for
192 the ESM-1b model⁴². We attempted to explain our model's unique ability to infer pseudo time with only the RBD
193 region. We explored the effectiveness of labels in our supervised learning, as it provides alternative directions
194 rather than the language model preference^{42, 44}. Interestingly, we found that the model prediction scores alone
195 have an even higher Spearman correlation score of 0.65 ($p < 1e-308$) with the sampling time (Fig. 3b) compared
196 with that of the inferred pseudo time, while for the predicted antibody escape potential, the Spearman correlation
197 is 0.67 ($p < 1e-308$). These findings verify our assumptions: under the immune selection pressure, the virus
198 evolves in the direction of immune escape, and our model can capture the antibody escape potential of the viral
199 variants.

200 We next assessed the antigenic evolution on a short time scale by comparing the model predictions against the
201 sampling time (Fig. 3c, Extended Data Fig. 4). We evaluated three types of scores, the ACE2 binding score, the

202 antibody escape potential, and the weighted average of the two scores. The predictiveness of the antibody escapes
203 score increases from nearly noninformative early in the pandemic to a stronger correlation during the Omicron
204 wave. It also gains predictiveness with the emergence and spread of Alpha variants in Early 2021 but
205 subsequently loses the predictiveness along with the emergence of other variants. We noted that the antigenic
206 evolution For the ACE2 binding probability score, it tends to become more informative during the first year,
207 while soon it becomes non-informative when the new VOC like Delta and Omicron emerged. These results
208 suggests that the antigenic evolution happens along with the infection waves.

209 We then examined the model sequence representations against the binding specificities. We found that after the
210 training, there are strong correlations between embeddings' primary and secondary axis of variation and the
211 binding specificities for all nine targets (Fig. 3d, Supplementary Fig. 5). The correlations are observed for both
212 ACE2 binding and antibody escape, suggesting that our multi-task learning strategies enable the model to learn
213 the functional properties simultaneously. Given that the variant sequence embeddings are shared across tasks,
214 this suggests that our model split the sequences based on an antigenic meaningful sense of binding preference.
215 In summary, our model effectively infers the immune escape potential and the ACE2 binding specificity, while
216 the predicted scores correlate positively with the real-world sampling time, especially for the newly emerging
217 Omicron wave. Taken together, we hypothesize that our model can work as a good scoring function for searching
218 for high-risk mutations and the corresponding variants.

219 ***In silico* directed evolution as a predictive tool**

220 With our model as the scoring function, we used the genetic algorithm to search for novel RBD variant sequences
221 with high antigenic evolutionary potential. The search process consists of selecting an initial sequence from the
222 GISAID database, generating and selecting “better-than-initial” sequences with the genetic algorithm to produce
223 38870 putatively high-risk variants within a 15-mutations “trust radius” of the initial sequence (Methods). We
224 performed the search process for the sequences in the GISAID database from January 1, 2022 to March 8, 2022,
225 yielding a total of 971 distinct sequences. We then visualized the generated sequences together with the existing
226 sequences using the distance-preserving multidimensional scaling plot⁴⁵ (Fig. 4a). While the sequences from the
227 deep mutational scanning experiments only occupy a small region around the wild type sequences, the prevalent
228 variants (e.g., Omicron) locate in different regions, far from the wild type. The sequences searched with our

229 model shown are diverse, largely expanding the sequence space.

230 Compared with the seed sequences, the synthetic sequences generated by our model include key mutations for
231 ACE2 binding and antibody escape. To visualize the difference and further explore the patterns of the generated
232 mutations, we constructed the position frequency matrix (PFM) for the two sequence sets and calculated the
233 Kullback-Leibler divergence (KL divergence) for each position based on the two PFMs (Methods). Fig. 4b and
234 Supplementary Fig. 6 provides structure-based visualizations and projects the Kullback-Leibler divergence per
235 site onto a crystal structure of the RBD (PDB id: 6m0j). As an alternative representation, Fig. 4c provides a
236 probability-weighted Kullback-Leibler logo plot⁴⁶ for the top 50 most divergence sites, where the total height of
237 the letters depicts the KL divergence of the site, while the size of the letters is proportional to the relative log-
238 odds score and observed probability (see Methods). The logo plot for all positions can be found in Extended
239 Data Fig. 5. Enriched amino acids locate at the positive side of the y-axis and depleted amino acids locate at the
240 negative side.

241 The logo plot shows that the mutations searched by our model largely overlap with the antibody escape maps.
242 For example, Y453, F456, and A475 are key sites for class 1 antibody escape¹⁸, while they are also present in
243 many synthetic variant sequences. Mutations escaped class 2 antibodies at sites E484, F490, and P491¹⁸. The
244 logo plot shows that these sites ranked high as “active sites”. Class 3 antibodies, which bind the opposite side of
245 the receptor-binding motif, tend to be escaped by sites like N437, N448, and Q498¹⁸, which are also vulnerable
246 sites suggested by the model. Class 4 antibodies bind to a conserved motif among the sarbecorviues, far away
247 from the RBM. Our model still captures the conservation and assigns mutations to the motif. However, some
248 sites with a large KL divergence do not locate in the epitope regions. This has several explanations. Firstly, it is
249 clear that some top sites (e.g., L368, C480) are not the direct binding sites but the proximal contact sites in the
250 structures, which may influence the binding as well. Secondly, as there are epistasis relationships among the
251 mutations, some combinatorial mutations may influence the RBD function nonlinearly and then modify the
252 antibody escape, which is not directly revealed by the epitope map. Moreover, these non-epitope sites with high
253 KL divergence need to be taken into consideration as they may perform an important role in future variants.
254 Another concern is that some sites in the epitope region have a low KL-divergence, one possible explanation is
255 that these sites have no tolerate mutations, for example, G416 and R457. Another explanation is that some

256 mutations at antibody-contact sites do not directly influence antibody binding.

257 The synthetic variant sequences share similar mutations with the chronic SARS-COV-2 infections. A reverse
258 mutation, R493Q, for example, was found in a persistently infected, immunocompromised individual⁴⁷. Other
259 mutations found by our model, like E340K⁴⁸, E484T³³, G485R⁴⁹, and F490L/E484G⁵⁰, are also found in
260 immunocompromised patients treated with monoclonal antibodies. Moreover, the unique mutations found in the
261 emerging variants, BA.4/5, the L452R, F486V, and the reverse mutation R493Q, are captured by our model. For
262 the newly emerging variants like XBB.1.5, the key mutation, F486P⁵¹, is also captured by our model (Extended
263 Data Fig. 5, Supplementary Fig. 7). This suggests that our model could be used for finding novel mutations that
264 may occur naturally. A detailed list of the found mutations in compromised patients is available in Supp Table
265 44. We next evaluated the immune evasion potential posed by the variant sequences using Evo-velocity analysis
266 and viral language model risk inference, followed by structure modeling and antibody-antigen docking. The
267 computational validation experiments suggests that the generated variants have high immune escape potential.
268 Further details can be found in the Supplementary Data 1-2.

269

270 **In vitro validation of novel mutations found by MLAEP**

271 Having generated the synthetic sequences and found interesting single mutations, it is thus crucial to validate the
272 risk and the immune evasion ability of combinatorial novel mutations using in vitro neutralizing antibody
273 binding assay, especially for those that cannot be predicted with a linear additive model. Though the Omicron
274 and its sub lineage are desired targets, they already exhibit high antibody escape abilities on the eight antibodies
275 we selected for training our model, making it difficult to distinguish the effectiveness of novel mutations induced
276 by MLAEP. To envision the differences, we used the RBD sequence of the Delta variant as the initial state and
277 ran the entire framework again to generate and select “better-than-Delta” sequences. Our goal was to find
278 possible antigenic evolutionary pathways for Delta that lead to high immune evasion.

279 We generated 3876 putatively high-risk variants using MLAEP and selected eight variants (Figure 5, Extended
280 Data Fig. 6) with unique immune evasion properties, including epistatic and non-epitope mutations. For example,
281 the RBD3 contains seven mutations compared to the wild type, but all the single mutations are experimentally
282 validated¹⁸ to be ineffective at evading the eight antibodies we used. However, our model predicted that the

283 RBD3 would have high immune evasion. The RBD4 does not contain mutations on the Class 4 antibody epitope,
284 but our model predicted that it would escape Class 4 antibodies. The selection criteria are detailed in
285 Supplementary Table 5.

286 We first expressed and purified the eight neutralizing monoclonal antibodies and ten RBDs (including wild type,
287 Delta, and eight synthetic RBD we generated) bearing different mutations. We tested different combinations of
288 neutralizing antibodies and RBD variants in a Homogeneous Time-Resolved Fluorescence (HTRF) based
289 antigen-antibody binding assay. In our HTRF-based binding assay, the wild type and Delta variant RBDs
290 exhibited high binding efficacy against different neutralizing monoclonal antibodies, with the IC₅₀ falling in
291 between 0.2 nM and 1 nM (Fig. 5). Notably, the Delta variant RBD showed no binding interaction to COV2-
292 2096 (Fig. 5), consistent with the literature that the L452R¹⁸ mutation on Delta variant confers evasion ability
293 against this neutralizing antibody. Intriguingly, all our predicted synthetic variants exhibited reduced or
294 diminished binding efficacy against all four classes of neutralizing antibodies targeting different epitope regions
295 (Fig. 5). Specifically, RBD4, RBD7, RBD8, and RBD9 exhibited evasion or reduced binding to COV2-2094 and
296 COV2-2677, two representative class 4 neutralizing monoclonal antibodies, even without bearing any mutations
297 in the class 4 epitope region. We also found that RBD8 could completely escape class 3 antibodies (COV2-2096
298 and COV2-2499) without bearing mutations in the class 3 epitope region, suggesting that epistasis relationship
299 play significant roles in the immune evasion, and such relationships could be captured by our deep learning
300 model. The RBD5, RBD7, and RBD8 variants retained sensitivity to class 1 (COV2-2832, COV2-2165) and
301 class 4 (COV2-2094, COV2-2677) antibodies with similar IC₅₀ values compared to wild type RBD, but their
302 binding efficacy to these neutralizing antibodies were reduced by large degrees. Overall, the synthetic variants
303 and the novel combinatorial mutations generated from MLEAP exhibited a high potency for immune evasion,
304 suggesting MLAEP captures the antigenic evolutionary potential.

305

306 **Discussion**

307 In this paper, we proposed a machine learning-guided antigenic evolution prediction paradigm for forecasting
308 the antigenic evolution of SARS-COV-2. We trained a multi-task deep learning model to predict ACE2/antibody
309 binding specificity using variant sequences and binding target structures. Predicting ACE2 binding specificity is

310 a relatively easy task, as one can capture the binding specificity using the unsupervised learning-based models²⁹.
311 However, predicting antibody binding specificity is much more challenging and less explored in the literature.
312 Through various validation experiments, we showed that our model can predict the antigenic evolutionary
313 potential resulting of high immune pressure. Combined with the genetic algorithm, we conducted *in silico*
314 directed evolution using the model scores. The resulting synthetic sequences displayed high immune evasion
315 potential, which we further validated using *in silico* computational tools and *in vitro* neutralizing antibody
316 binding assay. MLAEP captures mutations that also happen in chronic SARS-COV-2 infections and emerging
317 variants like BA.4/5 and XBB.1.5. In addition, MLAEP forecasts novel combinatorial mutations that affect
318 antibody binding beyond epitope regions. While we used the genetic algorithm to search for novel variants, other
319 search algorithms like hill-climbing⁵², simulated annealing⁵³, and reinforcement learning⁵⁴ could also be
320 combined with MLAEP. The multi-task learning model could be also replaced with other mutation effects
321 prediction models³⁰.

322 Deep learning models can learn high-order epistasis relationships among the multiple mutations^{28, 29}. Our multi-
323 task model, meanwhile, can capture such relationships and work as a monitor for predicting the escaping
324 potential of newly emerging variants, particularly heavily mutated variants. Our *in vitro* HTRF-based high
325 throughput assay verified that MLAEP is able to forecast epistatic and non-epitope mutations, thus expanding
326 our understanding and ability to predict the virus evolution. When combined with the Evo-velocity analysis, our
327 model helps to reveal the evolution trajectory of existing sequences and enables the discovery of high-risk
328 variants that may appear in the future. The results suggest that the *in silico* directed evolution can lead to the
329 prediction of *in vivo* virus evolution. Consequently, MLAEP may enable the support of public health decision-
330 making and guide the development of new vaccines. Besides, our approach could also be applied to the rapidly
331 evolving viruses and other potential outbreaks, such as antibiotic resistance⁵⁵.

332 An important property of MLAEP is that we focused on predicting the directionality of the mutation effect (i.e.,
333 whether a mutation increases or decreases binding affinity) rather than the magnitude of the effect. We plan to
334 further develop our model to capture the quantitative effect of mutations in the future. Besides, one limitation of
335 our model is that we only focused on the RBD sequences, while many mutations occur outside the region. We
336 noted that the mutagenesis assayed data provides semantically meaningful directions for finding “better-than-

337 natural” sequences. An increasing number of experiments characterize the functionality of mutations in other
338 regions, and we plan to explore these datasets in the future. Another concern is that we only optimized two
339 targets, the ACE2 binding and antibodies escape, while the directionality of evolution is also driven by many
340 other properties, like the epidemiology features and T cell responses. In addition, the limited availability of
341 variant ACE2 datasets prevented our model from capturing the fitness landscape of ACE2. Furthermore, the
342 virus evolves continuously, making the set of effective neutralization antibodies change over time. Fortunately,
343 the increasing availability of deep mutational scanning datasets^{19, 56} makes it convenient to track and update our
344 model regularly. In the future, we will use these datasets and incorporate more *in vivo* and *in vitro* experimental
345 data. Specifically, we will combine the *in vivo* antibody-antigen co-evolution data from patients and the
346 assessment of other immune responses to better understand and predict the evolution of SARS-COV-2.

347

348

349 **Method**

350 **Dataset**

351 We collected and cleaned nine deep mutational scanning datasets, which measure the binding affinity of the
352 SARS-COV-2 RBD variants towards the ACE2 and eight antibodies from four classes. We built a dataset
353 consisting of 19132 RBD sequences, where each sequence has nine labels, corresponding to their binding ability
354 to the nine targets. Most sequences have one or two mutations compared to the wild type RBD sequence.
355 However, considering the possible batch effect and the physical meaning differences among the measured scores,
356 we normalized each score independently by transforming the continuous variables into semantically meaningful
357 binary labels. For the ACE2 task, we directly compared the binding score of the mutated sequences to the wild
358 type and set the label to “enhanced binding” if the score is larger than the wild type and vice versa. There is no
359 information about the wild type for the eight antibodies tasks, so we cannot set the threshold as described earlier.
360 Instead, we found that the distributions of the binding score's logarithm clearly show two clusters; therefore, for
361 all antibody datasets, we took it as a mixture-of-Gaussian model respectively, defining the one with smaller
362 binding scores as escaped and vice versa. This preprocessing step is consistent with the subsequent work⁵⁷.
363 In summary, there were 1540 (8%) mutated RBD sequences identified as enhanced binding to ACE2, 3482 (18%)

364 mutated RBD sequences identified as escaped to COV2-2096, 1220 (6%) mutated RBD sequences identified as
365 escaped to COV2-2832, 2000 (10%) mutated RBD sequences identified as escaped to COV2-2094, 1473 (8%)
366 mutated RBD sequences identified as escaped to COV2-2050, 1859 (10%) mutated RBD sequences identified
367 as escaped to COV2-2677, 929 (5%) mutated RBD sequences identified as escaped to COV2-2479, 780 (4%)
368 mutated RBD sequences identified as escaped to COV2-2165 and 3347 (17%) mutated RBD sequences
369 identified as escaped to COV2-2499.

370 We used the pseudovirus neutralization test assay data from Liu et al.⁴¹ to validate our model performance. The
371 dataset measures the immune escaping of 10 high risk variants pseudoviruses by comparing the fold change in
372 IC_{50} of 17 monoclonal neutralizing antibody response against wild type pseudovirus.

373

374

375 **Overview of the multi-task model**

376 A central feature of SARS-COV-2 is antigenic evolution, that is, under high immune pressure, the newly
377 emerging variants will tend to escape the antibody while do not lose much binding ability to the ACE2.

378 To accomplish the goal of predicting the antigenic evolution, we need to construct the virtual fitness landscape
379 of the antigenic regions, especially for the RBD protein. We aimed to infer the fitness landscape of the RBD by
380 learning the effects of mutations on ACE2 binding and antibody escape. Specifically, given RBD variant
381 sequences and their labels, together with the binding partner (ACE2/antibody) structures, our model learned the
382 nonlinear mapping function f that can simultaneously predict the binding specificity for ACE2 and antibody. The
383 function f is parameterized by learnable mapping parameters θ composed of three modules: the sequence
384 feature extractor \mathcal{S} , the structure feature extractor \mathcal{G} and the sets of nine classification heads $H = \{\mathcal{H}_c\}_{c=1}^9$,
385 where all \mathcal{H}_c share the same group of parameters. All three modules are neural networks. The parameters of
386 the three modules are optimized in an end-to-end manner.

387 1. Sequence feature extractor. The sequence feature extractor takes as input of amino acid sequences
388 of RBD variants $x = (x_1, x_2, \dots, x_l)$ of length L , where L denotes the length of the RBD
389 sequence and elements x_i belongs to $A = \{\text{all amino acids}\}$. Input is mapped to a dense
390 representation vector (sequence representation). The backbone of the sequences feature extractor

391 is the ESM-1b transformer, which is pretrained on UniRef50 representative sequences with the
392 masked language modeling objective. We chose the ESM-1b as the sequence feature extractor
393 because it outperforms other baselines on a range of downstream tasks. The pretrained weights
394 were used for initializing the neural network, and we fine-tuned the model parameters during
395 training.

396 2. Representing structure as graph. We first represented the 3D structure as a k-nearest neighbor graph
397 $g = (V, E)$ with the node set $V = \{v_j\}_{j=1}^N$ of size N , where each element v_j denotes for the
398 features of representative atoms (we chose N, C, and O atom in the experiment) in the protein 3D
399 structure, N denotes the total number of atoms. For each atom, we got its two nearest neighbors
400 with the following constraints: the ones with the same atom type but belong to different amino
401 acids. We then measured the dihedral angles of the atom and its neighbors to as node features. The
402 edge features $E = \{e_{ij}\}_{i \neq j}$ describes the relationship between the nodes, including the relative
403 distance, direction, and orientation between the two nodes in the three-dimensional space. We set
404 k as 30.

405 3. Structure feature extractor. The structure feature extractor tasks as input of graphs $g = (V, E)$
406 describing the spatial feature of the protein structure. The transformed graph is further mapped into
407 a dense representation vector (structure representation). The backbone of the structure feature
408 extractor is a Structured Transformer³⁷, where the attention for each node is restricted to its k -
409 nearest neighbors in 3D space. We chose the Structured Transformer for the structure feature
410 extraction as it is computationally efficient and performs well in the protein design task. The
411 structure representation works as conditional tags in our multi-task learning.

412 4. Classification heads. After getting the sequence representation and the structure representation, we
413 concatenated the two vectors into the joint representation, and fed it into the classification heads.
414 The classification heads map the joint representation to the labels. We used nine parallel

415 classification heads for the nine classification tasks, while the neural network parameters are shared.
416 During training, the sequence feature extractor, the structure feature extractor and the classification
417 heads are trained in an end-to-end manner to minimize the average classification loss among the
418 nine tasks.

419 5. Loss function. Let $\mathbf{x} = \{x_i\}_{i=1}^N$ be the set of RBD variant amino acid sequences, and $\mathbf{y} =$
420 $\{y_i\}_{i=1}^N$ be the set of labels of all sequences, and $\mathbf{y} = \{y_i\}_{i=1}^N$ denotes for the set of M labels of
421 the i -th RBD variant. Furthermore, let $G = \{g_c(V, E)\}_{c=1}^M$ consists of M graphs derived from the
422 ACE/antibody structures. We seek to learn a joint embedding for all downstream classification
423 tasks to better model the fitness landscape of RBD. Therefore, the sequence feature extractor and
424 the structure feature extractor are shared among all tasks. Considering that all the tasks are
425 imbalanced in terms of the positive and negative samples, we added a rescaling weight p_c to all
426 tasks and optimized the following loss function:

$$427 \quad L = \frac{1}{MN} \sum_{c=1}^M \sum_{i=1}^N -[p_c y_i^c \cdot \log \sigma(\mathcal{H}_c(\mathcal{S}(x_i); \mathcal{G}(g_c))) + (1 - y_i^c) \cdot \log(1$$

428 $\quad \quad \quad - \sigma(\mathcal{H}_c(\mathcal{S}(x_i); \mathcal{G}(g_c))))]$

429
430 Where p_c equals to the number of positive samples divided by the number of negative samples,
431 M equals to nine, σ is the sigmoid function. The equation measures the binary cross entropy
432 between the targets and predicted probabilities.

433 434 **Architecture and hyperparameters**

435 The architecture of the sequences feature extractor is based on the ESM-1b transformer, which consists of 34
436 layers, we used the outputs of the 33rd layer as the sequence feature representations. For the structured
437 transformer, we only kept the transformer encoder, and used three layers of self-attention and position-wise

438 feedforward modules with a hidden dimension of 128. Finally, we got a 1280-d vector for each sequence as the
439 sequence representation and a 1300-d vector for each 3D structure as the structure representation. For each
440 classification head, we used 1024 neurons in the first layer and two neurons in the second layer. The RELU
441 function is used between the layers as nonlinear activations. We also passed a dropout rate $p=0.5$ and added
442 weight decay to prevent overfitting. We trained the entire model with the AdamW optimizer and used a linear
443 schedule with warmup to adjust the learning rate. We set the batch size as 16 and gradient accumulation steps as
444 10, which means that the total train batch size is 160, and the validation is the same. We used a weighted random
445 sampler function for our training batches, which oversamples the minority class to ensure that the number of
446 samples in each class are equal or close to equal. The model was trained for 9500 updates with the initial learning
447 rate of $1e-5$ and warmup steps 120, during which the model with the best marco F1 score among all the tasks
448 was kept. The hyperparameters described above were decided through several trials of experiments and selected
449 the one with the best performance.

450

451 **Choice of baselines**

452 Our framework follows the machine learning-guided directed evolution paradigm, while the quality of generated
453 sequences largely depends on the sequence-function model. Thus, we validated the generalization ability of the
454 models to newly seen variants with cross-validation. Our multi-task deep learning model was designed as a
455 supervised technique to infer the fitness of variant sequences, while there is no existing method designed for
456 multi-task learning or multi-label learning on this task. So instead, we evaluated the performance of the existing
457 methods separately for all nine tasks. We chose the state-of-the-art supervised-learning-based methods for
458 inferring the effect of mutations, including the augmented Potts model³⁰, the eUniRep model²⁶ and the gUniRep
459 model³⁸. We also benchmarked several baseline machine learning methods including CNN, LSTM, RNN, Linear
460 Regression, SVM, and Random Forest.

461 The augmented Potts model combines the evolutionary information with the one-hot encoded amino acid
462 sequences as input features and trains a linear regression model on top of the features. It outperforms most
463 existing methods in inferring the effects of mutations. We first generated the multiple sequence alignments
464 profile of RBD using the profile HMM homology search tool Jackhmmer. We set the bit score threshold as 0.5

465 and the number of iterations to 1. We then calculated the evolutionary Potts potential of the RBD variant
466 sequences using the plmc. We replaced the Ridge regression with a Logistic regression head for the classification
467 objective while keeping the rest procedures the same as the original settings.

468 The gUniRep model was trained on 24 million UniRef50 amino acid sequences with the next amino-acid
469 prediction objective, and the representations extracted from the pretrained model acts as a featurization of the
470 sequences, benefits the downstream protein informatics tasks. With the RBD variant sequences as input, we got
471 the fixed-length vector representations from the pretrained model as sequence embeddings. We added a Logistic
472 regression head for downstream classification.

473 The eUniRep model was built on top of the gUniRep model. An unsupervised fine-tuning step with sequences
474 related to the target protein (evotuning) was introduced to learn the distinct features of the target family. Previous
475 *in vitro* studies on the GFP and beta-lactamase proved its effectiveness for efficiently modeling the protein fitness
476 landscape. We performed evotuning with the same MSA profile we generated in the augmented Potts model.
477 After that, we characterized the sequence embeddings with the eUniRep model and train a Logistic regression
478 model for downstream classification. All methods were trained and tested on the same training data and
479 validation data for all five folds. The

480 The CNN model consisted of a feature module and a classification module, where the feature module was
481 composed of two 1D convolution layers, max-pooling layers, and ReLU layers. A 128-dimensional feature vector
482 generated by the feature module was used to predict the label by the classification layer.

483 The LSTM model consisted of a feature module and a classification module, where the feature module was
484 composed of one 1D LSTM layers and two linear layers, followed a ReLU layer and a sigmoid layer. A 128-
485 dimensional feature vector generated by the feature module was used to predict the label by the classification
486 layer. The RNN model is similar to LSTM model, except replace the LSTM module with the RNN module. The
487 Linear regression, SVM and Random Forest were implemented using Scikit-learn v1.1.0⁵⁸.

488

489

490

491 **Ablation study**

492 We performed ablation studies to show the effectiveness of each module. We first explored the effectiveness of
493 the graphical representations. We replaced the structure features with the random Uniform noise $X \sim U(0, 1)$
494 and performed the multi-task learning with the same training details and procedures. Besides, we also explored
495 the importance of fine tuning by freezing the parameters of the ESM-1b model and only optimized the parameters
496 of the classification head.

497

498 **Performance evaluation**

499 We evaluated the multi-task learning model with the 5-fold cross-validation. We randomly split the dataset into
500 five folds. Each time, we used four folds as the training data and held out the remaining fold for validation. We
501 used the Accuracy, Precision, Recall, and F1 score to evaluate the classification performance across all models.
502 As all the nine classification tasks are imbalanced, we used the Macro-precision, Macro-recall, and Macro-F1-
503 score to get an unbiased evaluation. True positive (TP), true negatives (TN), false positives (FP), and false
504 negatives (FN) were measured by comparison between the prediction results produced by the model and the
505 ground truth in the validation set.

506

$$507 \quad Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$508 \quad Recall_c = \frac{TP}{TP + FN}, \quad Recall_{macro} = \frac{Recall_1 + Recall_2}{2}$$

$$509 \quad Precision_c = \frac{TP}{TP + FP}, \quad Precision_{macro} = \frac{Precision_1 + Precision_2}{2}$$

$$510 \quad F1_{macro} = 2 \times \frac{Precision_{macro} \times Recall_{macro}}{Precision_{macro} + Recall_{macro}}$$

511

512

513 **Generate virtual RBD variant sequences with the genetic algorithm**

514 To forecast the variants that follows the antigenic evolutionary potential, we applied the genetic algorithm for
515 searching the peaks of the fitness landscape described by our model. Inspired by Darwin's theory of natural
516 evolution, the genetic algorithm mimics the evolutionary process in the genome, where mutations, crossover and

517 selection happen, letting candidate solutions of a population with higher fitness scores have a higher probability
518 of surviving and producing the next generation of offspring. For SARS-CoV-2, it has been proven that similar
519 progress happens in immunocompromised infected patients who got treated with the monocle antibodies³³.
520 Hence, we used the genetic algorithm to model the antigenic evolution process and search for the potential risky
521 variants that might appear in the future. The genetic algorithm we used consists of the following steps:

522 1. Selection of initial sequences. For the Omicron related experiments, the initial input sequences
523 were obtained from the March 8, 2022 GISAID release⁹. We selected the RBD sequences from the
524 recent two months (i.e. from January 1, 2022 to March 8, 2022), resulting in 957 distinct RBD
525 sequences. For the Delta experiments, the initial RBD sequence is the Delta variant RBD sequence.
526 For each sequence, we created a generation P_0 of size S by perturbing the sequences S times to
527 generate a set of distinct modifications to the original sequence.

528 2. Perturb operation. For a given sequence $X = (x_1, x_2, \dots, x_n)$, we first randomly selected an amino
529 acid x_i , and got the K nearest neighbors of the selected amino acid according to the BLUSUM62
530 matrix. Secondly, we computed the fitness value when x_i is replaced with its neighbors, while
531 keeping the remaining set of words unchanged. We then picked the mutation with a probability
532 proportional to its fitness value. Finally, the selected mutated amino acid replaced the original one,
533 we got a new sequence. We set $K=20$ in our experiments.

534 3. Estimation of the fitness. The fitness score is defined as the average value of target label prediction
535 probabilities for all nine tasks. For ACE2 binding task, the target label is binding, while for the
536 antibody task, the target label is escape. The probabilities were found by querying the trained multi-
537 task model.

538 4. Crossover. After getting the perturbed population and the fitness values for each individual
539 sequence, we performed the crossover operation. Pairs of sequences are randomly selected with the
540 probabilities proportional to its fitness value. A child sequence is then generated by independently
541 sampling from the two parents. The newly generated sequences form the new generation. If the

542 fitness value of a population member in the generation is higher than the high-risk threshold, the
543 optimization is done. Otherwise, the perturbation, selection and crossover operation will be applied
544 on the new generation.

545 We performed *in silico* evolution for each initial sequence from GISAID subset for 100 times independently, and
546 finally got 38870 unique RBD variant sequences for the Omicron experiments. We got 3876 unique RBD variant
547 sequences for the Delta experiments. We selected eight RBDs for the HTRF-based neutralizing antibody binding
548 assay (Supp Table 5).

549

550 **Evo-velocity analysis**

551 The Evo-velocity analysis follows the study of Hie et al⁴². They used the pretrained protein language model (e.g.
552 ESM-1b) to predict the local evolution within protein families and used a dynamic “vector field” to visualize it.
553 It involves embedding the sequences of interest as vectors in a high-dimensional latent space, where the
554 geometric distance between the representation of proteins correlates with their actual structural, functional, and
555 evolutionary relatedness. The evo-velocity between two sequences is calculated by considering the log-
556 pseudolikelihood of observing a mutation from one sequence to another, providing a local mutational likelihood
557 gradient around a particular protein. When looked at globally, this vector field gives insight into the directionality
558 of the evolutionary process and can model global evolution. We first computed the embeddings for each sequence
559 with the ESM-1b model, and then constructed the K-nearest-neighbor graph based on the embeddings, in which
560 node represents the sequences and edges connected similar sequences. Further, the edges were assigned with
561 directions based on the language model pseudolikelihoods, with flow-in node meaning evolutionarily favorable.
562 Here, we performed the Evo-velocity analysis with their settings and ours. In our setting, we used the joint
563 embeddings extracted from the fine-tuned protein language model and the Structured Transformer model to
564 represent the sequences and set the direction of the edge by comparing the average predicted score among the
565 nine tasks, where vertex with a large value is defined as the tail. We collected 7594 unique RBD sequences from
566 the March 8, 2022 GISAID release. The date of the sequence is defined as the first reported date. After
567 constructing the directed KNN neighborhood graph, we further performed network diffusion analysis to infer

568 the pseudo time. We manually set the root as the wild type RBD sequences.

569

570 **Visualization**

571 We visualized the model embeddings using UMAP. The K-nearest-neighbor network was built with the k set to
572 30, while the resolution is set to 1. We calculated the KNN graph and performed UMAP for both the GISAID
573 sequences and the generated sequences. We further projected the predicted KL-divergence maps onto a RBD
574 structure (PDB id: 6m0j) and visualized the structure with PyMol. We collected the binding epitopes for class 1,
575 2, and 3 antibodies from Greaney et al.¹⁸, while for the class 4 antibodies, we used the contact sites of antibody
576 CR3022 to represent the class 4 binding epitope.

577 We used probability weighted Kullback-Leibler Logo plot for visualizing the generated mutations. Let $M_1 =$
578 $(f_1, f_2, f_3, \dots, f_n)$ denote the position frequency matrix (PFM) of the initial sequences, where the length of
579 the initial sequences are n and each $f_i = (a_1, a_2, \dots, a_{20})^T$, represents the frequency of each amino acid at
580 position i . Further, let $M_2 = (f'_1, f'_2, \dots, f'_n)$ denotes for the PFM of the generated sequences, each $f'_i =$
581 $(a'_1, a'_2, \dots, a'_{20})^T$. We computed the KL divergence for each position:

$$582 \quad D_{KL}(f'_i || f_i) = \sum_{i=1}^{20} a'_i \cdot \ln(a'_i/a_i)$$

583

584 The KL divergence denotes for the total heights at each position in the logo plot. We further set the height and
585 the direction of a letter with a probability weighted normalization⁴⁶, where the relative height of each individual
586 amino acid is proportional to $a'_i \cdot \ln(a'_i/a_i)$:

$$587 \quad h(a'_i) = \frac{a'_i \cdot \ln(a'_i/a_i)}{\sum_{i=1}^{20} a'_i \cdot |\ln(a'_i/a_i)|} D_{KL}(f'_i || f_i)$$

588

589 **Recombinant monoclonal antibody and RBD variants purification**

590 The sequences coding SARS-COV-2 monoclonal antibodies were kindly provided by Prof. James E. Crowe from
591 Vanderbilt University Medical Center. The LH and HC sequences were codon optimized and submitted to
592 Genescript for custom human IgG1 antibody expression. Sequences of wild type, deltavariant, and synthetic

593 variant RBD proteins were codon optimized and submitted to Twist for vector construction. All RBD constructs
594 contain a secretion signal on the N-terminal, and a 6× his tag followed by a strep-tag II on the C-terminal. In
595 brief, Expi293 cells were transfected in 40 mL Expi293 Expression Medium (Thermo Fisher A1435101) at 37°C,
596 8% CO₂ on an orbital shaker at 120 rpm. After five days, cells were removed by spinning at 500 ×g for 5 mins
597 at 4 °C, and the medium was further centrifuged at 16000 ×g for 5 mins at 4 °C. The supernatant was then mixed
598 with his-tag purification resin (Beyotime P2221) on a shaker at 4°C. After 1 hour of incubation, the mixture was
599 loaded on a gravity chromatography column and washed for 15 mL of washing buffer [25 mM Tris, pH 8, 300
600 mM NaCl, and 1 mM DTT]. The elution was collected in 5 mL and loaded on another 2 mL column pre-packed
601 with 0.5 mL Strep-Tactin XT 4Flow high-capacity resin (IBA Lifesciences 2-5030-025). The RBD proteins were
602 eluted in 5 mL of washing buffer supplemented with 50 mM Biotin. For some mutant RBD proteins that have
603 reduced secretion into the medium, cell lysates were prepared in lysis buffer [25 mM Tris, pH 8, 300 mM NaCl,
604 0.5 % Triton X-100, 1 mM DTT, 1× protease inhibitor cocktail (PIC)] for 30 min on a shaker at 4°C. Clarified
605 lysates were subject to two affinity columns following the same purification protocols. All purified RBD proteins
606 were buffer exchanged and concentrated to 1 μM in 1× PBS using Amicon, flash-frozen in liquid nitrogen, and
607 stored at -80 °C.

608

609 **Homogeneous Time Resolved Fluorescence (HTRF) antigen-antibody binding assay**

610 The binding intensity between purified SARS-COV-2 RBDs and neutralizing antibodies was measured as HTRF
611 signals in the antigen-antibody binding assay. The HTRF donor and acceptor pair was chosen to target the his-
612 tagged RBD proteins and human IgG1 antibodies, respectively. Briefly, a total of 10 μL reaction was set up on
613 each well of the black, round-bottom, low-volume 384-well plates (Corning 4511) containing 5 nM purified wild
614 type or mutant RBDs, 3 nM goat anti-human IgG conjugated with Alex Fluor 647 (Thermo Fisher A-21445),
615 0.33nM monoclonal antibody anti-6His-Tb-cryptate Gold (Cisbio 61HI2TLA) and two-fold dilutions of
616 neutralizing mAbs from 2 nM to 0.0156 nM in 1× PBS supplemented with 0.1 % BSA, and 0.1 % Tween-20.
617 The plate was sealed with plastic film and incubated at room temperature for 1 hour. The HTRF signals were
618 measured in CLARIOstar Plus (BMG LABTECH) with the excitation filter at 340 nm and the emission filters
619 at 620 nm and 665 nm. The reading lag time and integration time were set to 60 μs and 200 μs, respectively. The

620 HTRF ratios from samples and negative controls were calculated by dividing the intensity readouts from the 665
621 nm channel over the 620 nm channel. All ratios were background subtracted and normalized in ΔF %:

622

$$623 \quad \Delta F\% = \frac{\text{HTRF ratio(sample)} - \text{HTRF ratio(negative control)}}{\text{HTRF ratio(negative control)}} \times 100$$

624

625 The IC50 value was calculated by fitting the data into a dose-response curve in Prism 9. Data points with the
626 ‘hook’ effect were removed from the fitting.

627

628

629 **Data availability**

630 The deep mutational scanning datasets is publicly available at [https://jbloomlab.github.io/SARS-CoV-2-](https://jbloomlab.github.io/SARS-CoV-2-RBD_DMS/)
631 [RBD_DMS/](https://jbloomlab.github.io/SARS-CoV-2-RBD_DMS/) and https://github.com/jbloomlab/SARS-CoV-2-RBD_MAP_Crowe_antibodies. The
632 pseudovirus neutralization test assay data is publicly available in its original publications. The GISAID data is
633 publicly available at <https://www.gisaid.org/>. We retrieved data from the website on 8 March 2022.

634

635

636 **Acknowledgment**

637 We acknowledge all members in the Structural and Functional Bioinformatics (SFB) group at KAUST for the
638 fruitful discussions.

639

640 **Funding**

641 King Abdullah University of Science and Technology (KAUST). [FCC/1/1976-44-01, FCC/1/1976-45-01,
642 URF/1/4663-01-01, REI/1/5202-01-01, REI/1/4940-01-01, and RGC/3/4816-01-01] to X.G. University Grants
643 Committee’s Collaborative Research Fund (C6036-21GF) to P.P.H.C. The Chinese University of Hong Kong’s
644 Research Committee Research Fellowship to X.X.

645

646 **Author contributions**

647 W.H., C.N., X.X., P.P.H.C and X.G. conceptualized the study and developed the methodology. W.H. and C.N.
648 implemented models and analyzed data. X.X., P.P.H.C., R.Z., Y.W., and S.S. designed and performed the wet-
649 lab experiments. A.S. developed the web server. Z.L., H.Z., and J.Z. helped with the baseline experiments.
650 P.P.H.C and X.G. supervised the research and the entire project. All authors wrote the paper.

651

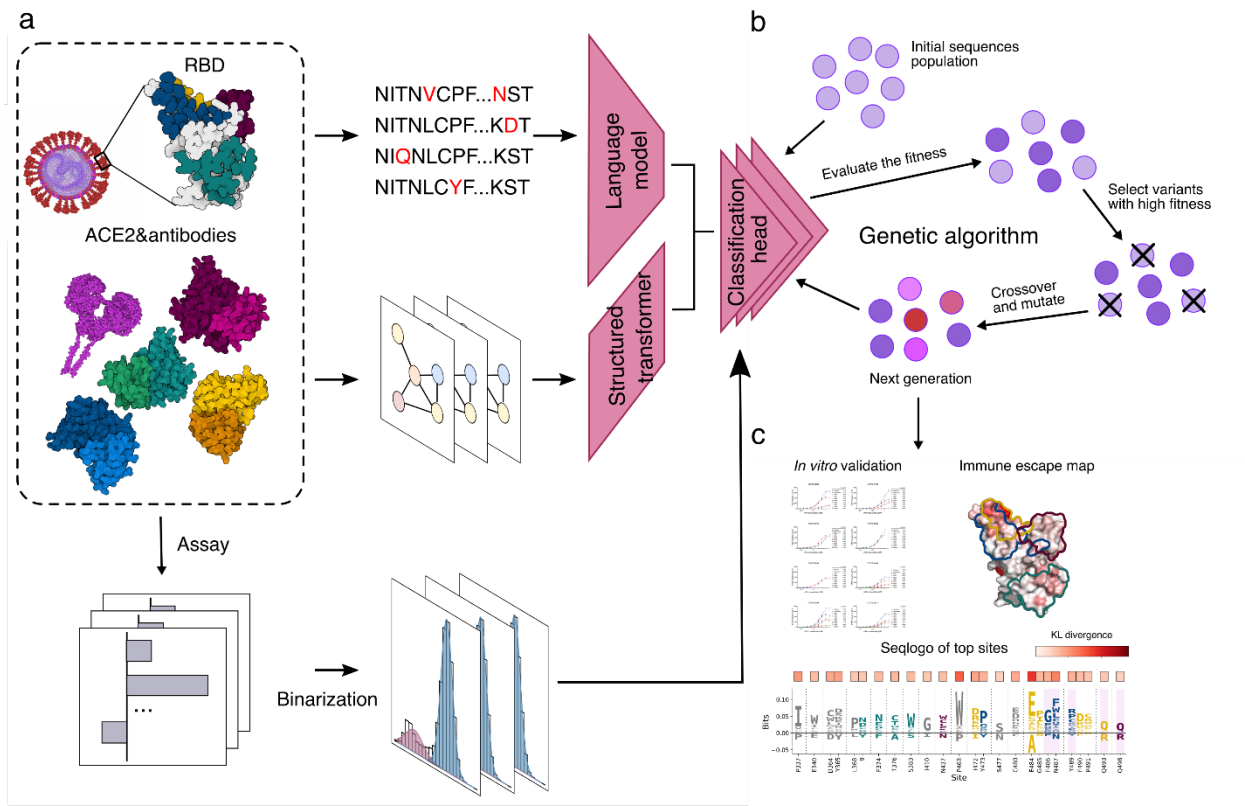
652 **Reference**

- 653 1. Wang, P. et al. Antibody resistance of SARS-CoV-2 variants B.1.351 and B.1.1.7. *Nature* **593**, 130-
654 135 (2021).
- 655 2. Wang, P. et al. Increased resistance of SARS-CoV-2 variant P.1 to antibody neutralization. *Cell*
656 *Host & Microbe* **29**, 747-751.e744 (2021).
- 657 3. McCormick, K.D., Jacobs, J.L. & Mellors, J.W. The emerging plasticity of SARS-CoV-2. *Science*
658 *(New York, N.Y.)* **371**, 1306-1308 (2021).
- 659 4. Iketani, S. et al. Antibody evasion properties of SARS-CoV-2 Omicron sublineages. *Nature* **604**,
660 553-556 (2022).
- 661 5. Starr, T.N. et al. Deep mutational scanning of SARS-CoV-2 receptor binding domain reveals
662 constraints on folding and ACE2 binding. **182**, 1295-1310. e1220 (2020).
- 663 6. Ozono, S. et al. SARS-CoV-2 D614G spike mutation increases entry efficiency with enhanced
664 ACE2-binding affinity. **12**, 1-9 (2021).
- 665 7. Jangra, S. et al. SARS-CoV-2 spike E484K mutation reduces antibody neutralisation. **2**, e283-
666 e284 (2021).
- 667 8. Volz, E. et al. Assessing transmissibility of SARS-CoV-2 lineage B. 1.1. 7 in England. **593**, 266-
668 269 (2021).
- 669 9. Charmet, T. et al. Impact of original, B.1.1.7, and B.1.351/P.1 SARS-CoV-2 lineages on vaccine
670 effectiveness of two doses of COVID-19 mRNA vaccines: Results from a nationwide case-control
671 study in France. *The Lancet Regional Health - Europe* **8**, 100171 (2021).
- 672 10. Mlcochova, P. et al. SARS-CoV-2 B.1.617.2 Delta variant replication and immune evasion. *Nature*
673 **599**, 114-119 (2021).
- 674 11. Ong, S.W.X. et al. Clinical and virological features of SARS-CoV-2 variants of concern: a
675 retrospective cohort study comparing B.1.1.7 (Alpha), B.1.315 (Beta), and B.1.617.2 (Delta).
676 *Clinical infectious diseases : an official publication of the Infectious Diseases Society of America*
677 (2021).
- 678 12. Andrews, N. et al. Covid-19 vaccine effectiveness against the Omicron (B. 1.1. 529) variant.
679 (2022).
- 680 13. Piccoli, L. et al. Mapping neutralizing and immunodominant sites on the SARS-CoV-2 spike
681 receptor-binding domain by structure-guided high-resolution serology. **183**, 1024-1042. e1021
682 (2020).
- 683 14. Dejnirattisai, W. et al. The antigenic anatomy of SARS-CoV-2 receptor binding domain. **184**,
684 2183-2200. e2122 (2021).
- 685 15. Zost, S.J. et al. Potently neutralizing and protective human antibodies against SARS-CoV-2.
686 **584**, 443-449 (2020).
- 687 16. Barnes, C.O. et al. SARS-CoV-2 neutralizing antibody structures inform therapeutic strategies.
688 **588**, 682-687 (2020).

- 689 17. Shrestha, L.B., Tedla, N. & Bull, R.A. Broadly-Neutralizing Antibodies Against Emerging SARS-
690 CoV-2 Variants. *Front Immunol* **12**, 752003-752003 (2021).
- 691 18. Greaney, A.J. et al. Mapping mutations to the SARS-CoV-2 RBD that escape binding by different
692 classes of antibodies. *Nature Communications* **12**, 4196 (2021).
- 693 19. Cao, Y. et al. Omicron escapes the majority of existing SARS-CoV-2 neutralizing antibodies.
694 *Nature* **602**, 657-663 (2022).
- 695 20. Hie, B., Zhong, E.D., Berger, B. & Bryson, B.J.S. Learning the language of viral evolution and
696 escape. **371**, 284-288 (2021).
- 697 21. Taft, J.M. et al. Predictive profiling of SARS-CoV-2 variants by deep mutational learning.
698 (2021).
- 699 22. Beguir, K. et al. Early Computational Detection of Potential High Risk SARS-CoV-2 Variants.
700 (2021).
- 701 23. Maher, M.C. et al. Predicting the mutational drivers of future SARS-CoV-2 variants of concern.
702 eabk3445 (2021).
- 703 24. Markov, P.V., Katzourakis, A. & Stilianakis, N.I. Antigenic evolution will lead to new SARS-CoV-2
704 variants with unpredictable severity. *Nature Reviews Microbiology* **20**, 251-252 (2022).
- 705 25. Starr, T.N. et al. Shifting mutational constraints in the SARS-CoV-2 receptor-binding domain
706 during viral evolution. (2022).
- 707 26. Biswas, S., Khimulya, G., Alley, E.C., Esvelt, K.M. & Church, G.M.J.N.m. Low-N protein
708 engineering with data-efficient deep learning. **18**, 389-396 (2021).
- 709 27. Yang, K.K., Wu, Z. & Arnold, F.H.J.N.m. Machine-learning-guided directed evolution for protein
710 engineering. **16**, 687-694 (2019).
- 711 28. Riesselman, A.J., Ingraham, J.B. & Marks, D.S. Deep generative models of genetic variation
712 capture the effects of mutations. *Nature Methods* **15**, 816-822 (2018).
- 713 29. Rives, A. et al. Biological structure and function emerge from scaling unsupervised learning to
714 250 million protein sequences. **118**, e2016239118 (2021).
- 715 30. Hsu, C., Nisonoff, H., Fannjiang, C. & Listgarten, J.J.N.b. Learning protein fitness models from
716 evolutionary and assay-labeled data. 1-9 (2022).
- 717 31. Whitley, D.J.S. & computing A genetic algorithm tutorial. **4**, 65-85 (1994).
- 718 32. Alzantot, M. et al. Generating natural language adversarial examples. (2018).
- 719 33. Halfmann, P. et al. Evolution of a globally unique SARS-CoV-2 Spike E484T monoclonal antibody
720 escape mutation in a persistently infected, immunocompromised individual. (2022).
- 721 34. Harari, S. et al. Drivers of adaptive evolution during chronic SARS-CoV-2 infections. (2022).
- 722 35. Wilkinson, S.A. et al. Recurrent SARS-CoV-2 mutations in immunodeficient patients. (2022).
- 723 36. UniProt: the universal protein knowledgebase in 2021 %J Nucleic acids research. **49**, D480-
724 D489 (2021).
- 725 37. Ingraham, J., Garg, V., Barzilay, R. & Jaakkola, T.J.A.i.N.I.P.S. Generative models for graph-based
726 protein design. **32** (2019).
- 727 38. Alley, E.C., Khimulya, G., Biswas, S., AlQuraishi, M. & Church, G.M.J.N.m. Unified rational protein
728 engineering with sequence-based deep representation learning. **16**, 1315-1322 (2019).
- 729 39. Hie, B.L. et al. Efficient evolution of human antibodies from general protein language models
730 and sequence information alone. 2022.2004.2010.487811 (2022).
- 731 40. Mason, D.M. et al. Optimization of therapeutic antibodies by predicting antigen specificity from
732 antibody sequence via deep learning. *Nature Biomedical Engineering* **5**, 600-612 (2021).
- 733 41. Liu, L. et al. Striking antibody evasion manifested by the Omicron variant of SARS-CoV-2. *Nature*
734 **602**, 676-681 (2022).
- 735 42. Hie, B.L., Yang, K.K. & Kim, P.S. Evolutionary velocity with protein language models predicts
736 evolutionary dynamics of diverse proteins. *Cell Systems* **13**, 274-285.e276 (2022).

- 737 43. McInnes, L., Healy, J. & Melville, J.J.a.p.a. Umap: Uniform manifold approximation and
738 projection for dimension reduction. (2018).
- 739 44. Sandhu, M., Spence, M.A. & Jackson, C.J. Evo-velocity: Protein language modeling accelerates
740 the study of evolution. *Cell Systems* **13**, 271-273 (2022).
- 741 45. Cox, M.A. & Cox, T.F. in Handbook of data visualization 315-347 (Springer, 2008).
- 742 46. Thomsen, M.C.F. & Nielsen, M. Seq2Logo: a method for construction and visualization of amino
743 acid binding motifs and sequence profiles including sequence weighting, pseudo counts and
744 two-sided representation of amino acid enrichment and depletion. *Nucleic Acids Research* **40**,
745 W281-W287 (2012).
- 746 47. Scherer, E.M. et al. SARS-CoV-2 evolution and immune escape in immunocompromised patients
747 treated with exogenous antibodies. 2022.2004.2012.22273675 (2022).
- 748 48. Destras, G., Bal, A., Simon, B., Lina, B. & Josset, L. Sotrovimab drives SARS-CoV-2 Omicron
749 variant evolution in immunocompromised patients. 2022.2004.2008.22273513 (2022).
- 750 49. Harari, S. et al. Drivers of adaptive evolution during chronic SARS-CoV-2 infections.
751 2022.2002.2017.22270829 (2022).
- 752 50. Wilkinson, S.A. et al. Recurrent SARS-CoV-2 Mutations in Immunodeficient Patients.
753 2022.2003.2002.22271697 (2022).
- 754 51. Yue, C. et al. Enhanced transmissibility of XBB.1.5 is contributed by both strong ACE2 binding
755 and antibody evasion. 2023.2001.2003.522427 (2023).
- 756 52. Davis, L. in Proc. Intl. Conf. Genetic Algorithm, 1991 18-23 (1991).
- 757 53. Van Laarhoven, P.J. & Aarts, E.H. in Simulated annealing: Theory and applications 7-15 (Springer,
758 1987).
- 759 54. Nareyek, A. in Metaheuristics: Computer decision-making 523-544 (Springer, 2003).
- 760 55. Li, Y. et al. Hmd-arg: hierarchical multi-task deep learning for annotating antibiotic resistance
761 genes. **9**, 1-12 (2021).
- 762 56. Dadonaite, B. et al. A pseudovirus system enables deep mutational scanning of the full SARS-
763 CoV-2 spike. 2022.2010.2013.512056 (2022).
- 764 57. Frazer, J. et al. Disease variant prediction with deep generative models of evolutionary data.
765 **599**, 91-95 (2021).
- 766 58. Pedregosa, F. et al. Scikit-learn: Machine learning in Python. **12**, 2825-2830 (2011).
- 767
- 768

769

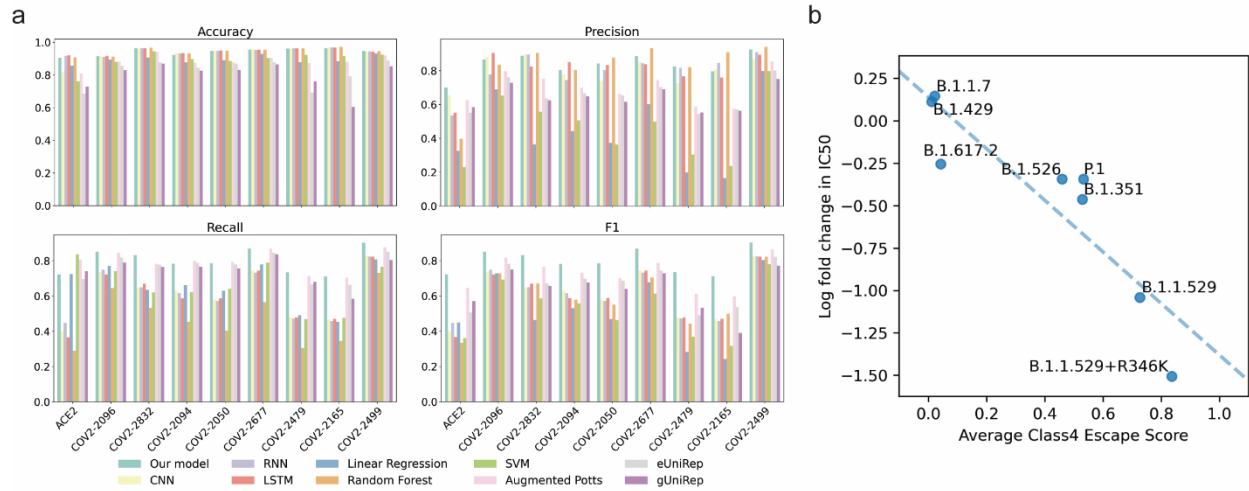


770

771 **Fig. 1 | Overview of the MLAEP framework.** a, The multi-task learning model. We collected and cleaned the
 772 RBD variant sequences and their corresponding binding/escaping specificity to the ACE2 and eight antibodies.
 773 Then, the sequences and the structures of their binding partners were fed into the deep learning model with the
 774 multi-task learning objective. b, The genetic algorithm. *In silico* directed evolution was performed to navigate
 775 the virtual fitness landscape defined by the nine scores from the multi-task model. The generation loop was
 776 repeated multiple times until the desired functionality was reached. c, These generated sequences were then
 777 subjected to validation experiments for evaluating their functional attributes.

778

779



780

781 **Fig. 2 | Performance evaluation and *in vitro* pVNT experimental data validation.** a, Model performance

782 comparison for the classification of ACE2 and antibody binding specificity across different algorithms. Including

783 our model, augmented Potts model, eUniRep model, gUniRep model, CNN, RNN, LSTM, linear regression,

784 SVM and random forest. The details of model implementation are given in Methods and performance metrics

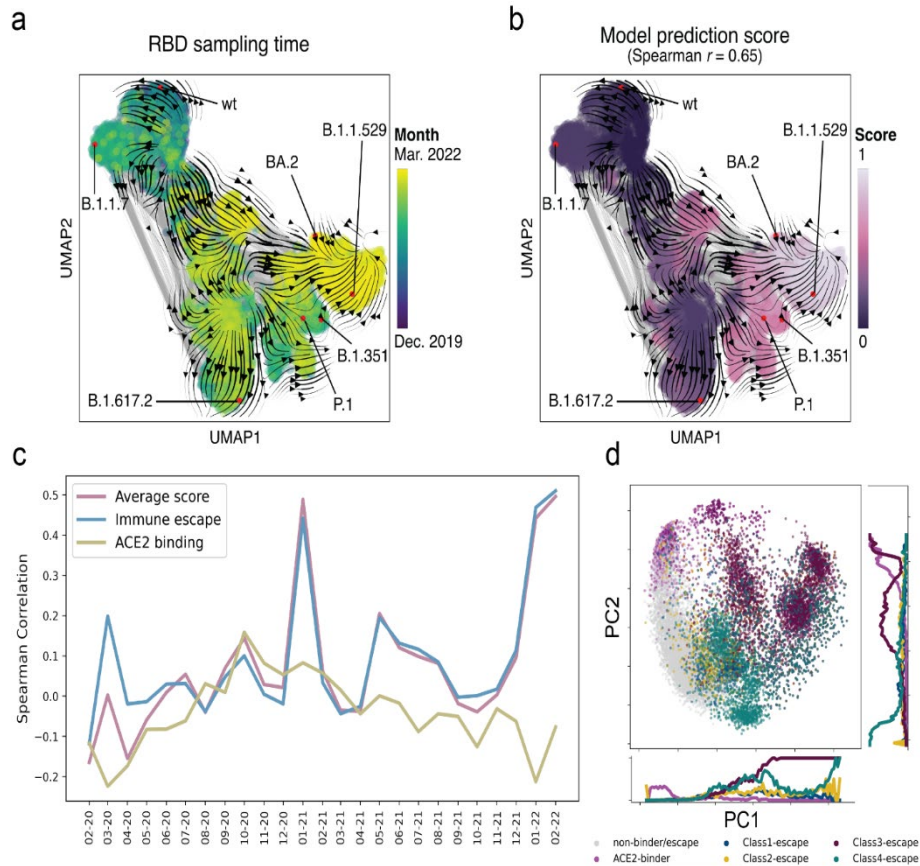
785 were calculated according to the equations provided in the Methods. b, Validation of the predicted immune

786 escape potential using the class4 monoclonal antibody-based pVNT assay data (Antibody 10-40). The x-axis

787 indicates the model predicted variant escape potential, while the y-axis is the log fold change of the VOCs

788 compared with the wild type.

789



790

791 **Fig. 3 | Multi-task model captures the antigenic evolutionary potential. a,** The landscape of SARS-COV-2

792 RBD variant sequences (obtained from GISAID), represented as a KNN-similarity graph (with the darker blue

793 region represents less recent date, e.g., 2019, and yellow represents more recent date, e.g., 2022). The gray lines

794 indicate graph edges, while the colored points are sequences with the known sampling time. The streamlines

795 among the points show a visual correlation between model predicted scores and the known sampling time. b,

796 Use the average score of our model to visualize the landscape. The landscape is colored by the model prediction

797 score with darker colors represent lower scores and lighter colors represent higher scores. c, Spearman

798 correlation overtime for the model predictions, including the ACE2 binding score, immune escape potential, and

799 the weighted average of the two in a time window of previous three months for each sampled date. (From

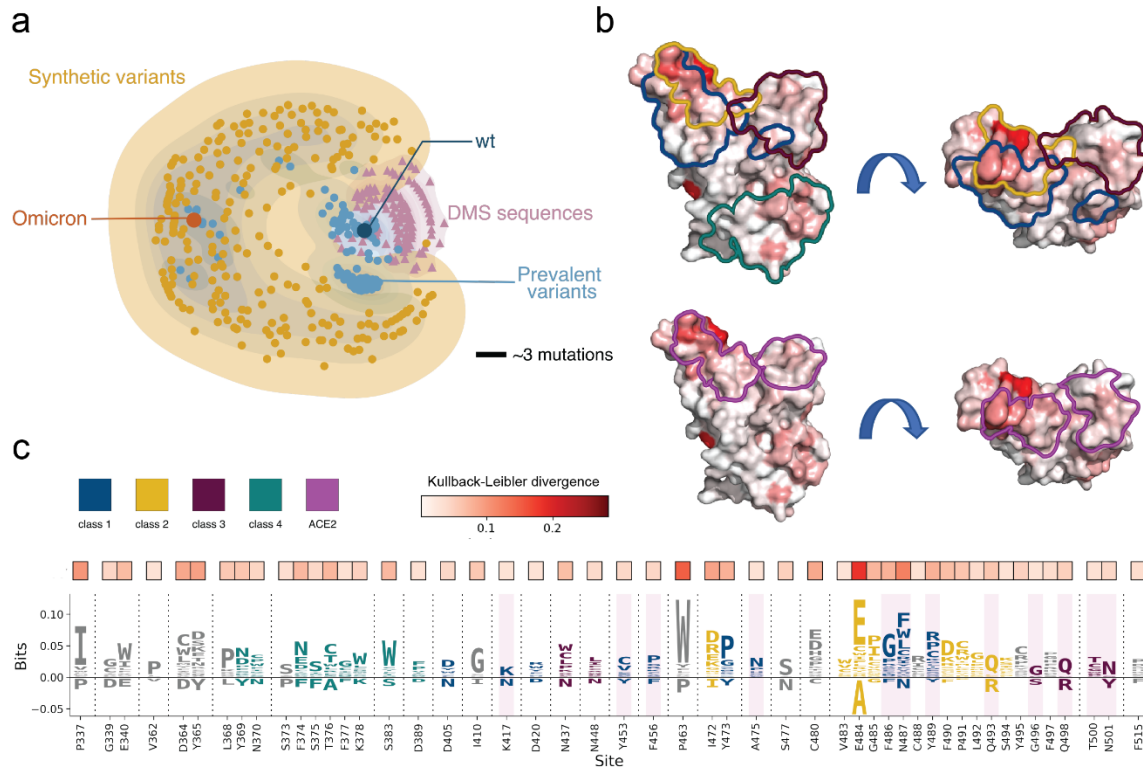
800 February 2020 to February 2022) d, Principal component analyses of the sequence's representations from our

801 model, colored by the escaping/binding ability towards COV2-2832, COV2-2165(class 1 antibody), COV2-2479,

802 COV2-2500 (class 2 antibody), COV2-2096, COV2-2499 (class 3 antibody), COV2-2677, COV2-2094 (class 4

803 antibody) and ACE2.

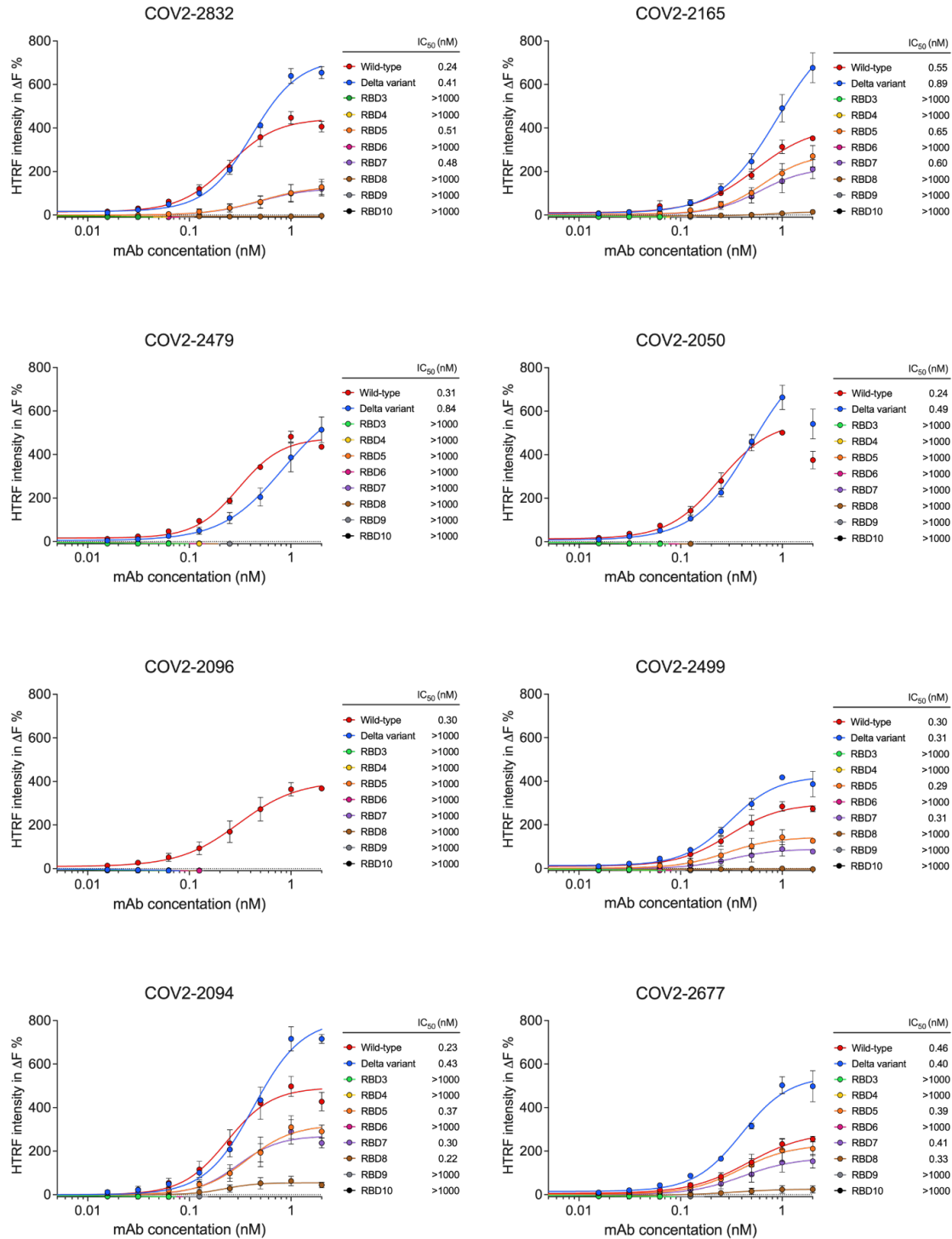
804



805

806 **Fig. 4 | Overview of the synthetic sequences.** a, Distance-preserving multidimensional scaling plot illustrates
 807 synthetic sequences' diversity compared to existing variants and deep mutagenesis sequences. A scale bar of
 808 three mutations is shown. b and c, the differences between the initial sequences and the synthetic sequences. b,
 809 The surface of the RBD protein, colored by the KL divergence between the initial sequences and the synthetic
 810 sequences. Colored outlines indicate the epitope structural footprint. c, The top 50 sites with the highest KL-
 811 divergence value are selected for visualizing the difference between the generated sequences and the existing
 812 sequences. Enriched amino acids locate at the positive side of the y-axis and depleted amino acids locate at the
 813 negative side.

814



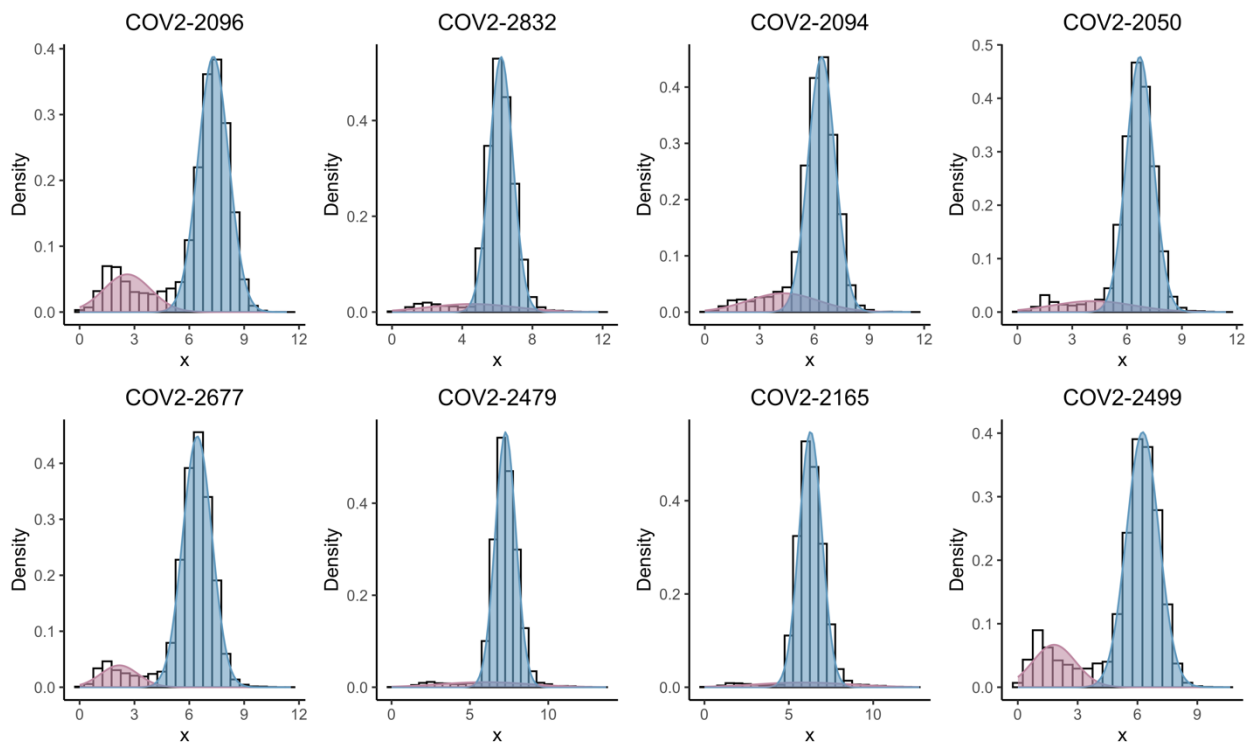
815
816
817
818
819
820
821
822

Figure 5. | Epitope mutations confer RBD resistance to the binding of neutralizing mAbs.

HTRF-based binding assay of wild type and mutant RBD proteins against two representative anti-RBD monoclonal antibodies from four classes, including COV2-2832 and COV2-2165 (class 1 antibody), COV2-2479 and COV2-2050 (class 2 antibody), COV2-2096 and COV2-2499 (class 3 antibody), as well as COV2-2094 and COV2-2677 (class 4 antibody). ΔF % values were calculated from raw data and fit into dose-response curves, and the IC_{50} values were listed side by side. Error bars represent standard deviation (n = 3).

823

824



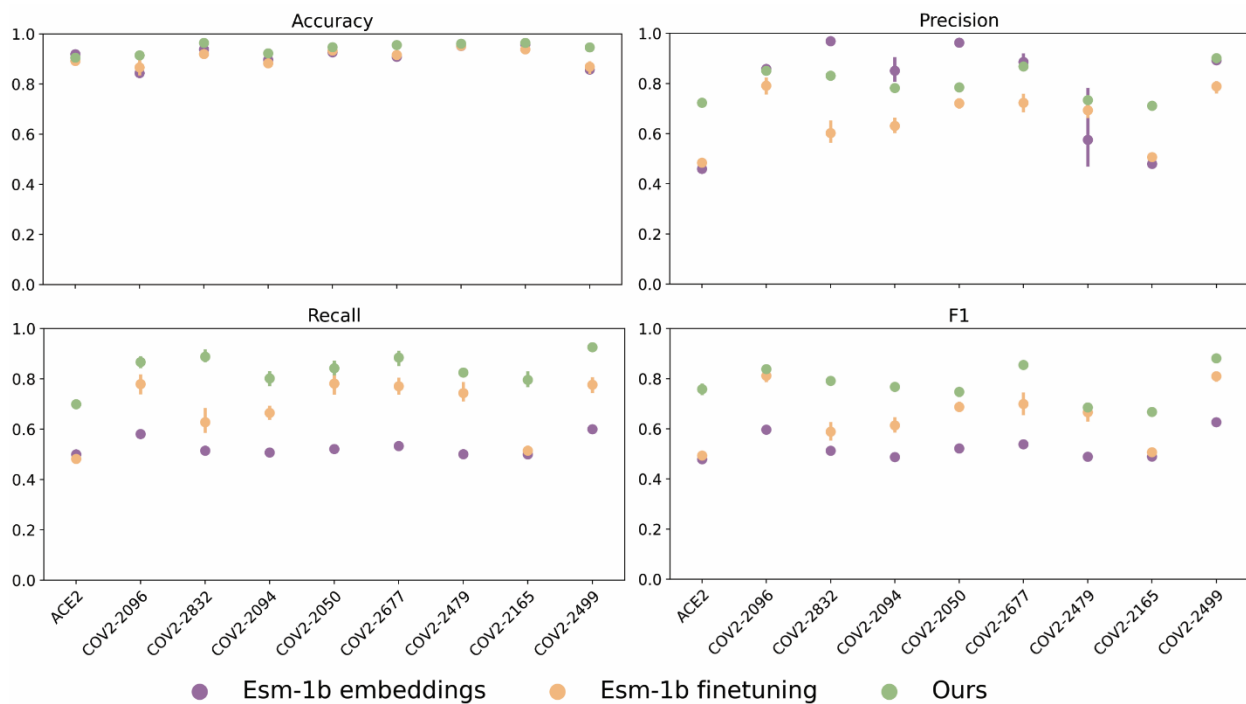
825

826 **Extended Data Fig. 1 | The distribution of the DMS scores of eight antibodies.** We log-transformed the deep
827 mutational scanning scores, and got clearly two clusters for all antibodies. We then used the Gaussian mixture
828 model to split the score into two clusters. Red cluster is defined as non-escape, while the blue clusters is defined
829 as escape.

830

831

832



833

834 **Extended Data Fig. 2 | Ablation studies.** Justification of using the fine-tuning, the structure representations in
835 the multi-task learning framework, in terms of Accuracy, macro-Precision, macro-Recall and macro-F1 score.
836 ESM-1b shows the results of the fine-tuning steps' ablation. ESM-1b finetuning shows the results of replacing
837 the Structured Transformer's ablation.

838

839

840

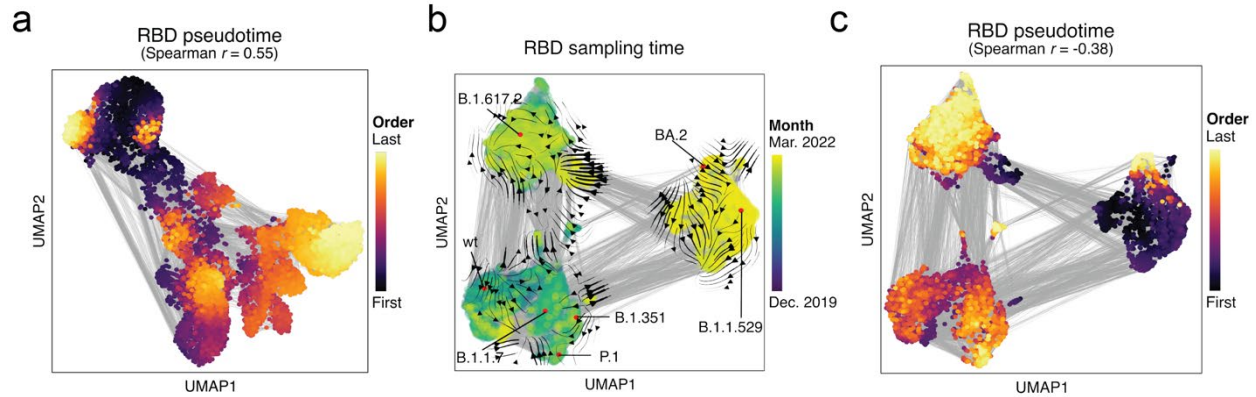
841

842

843

844

845

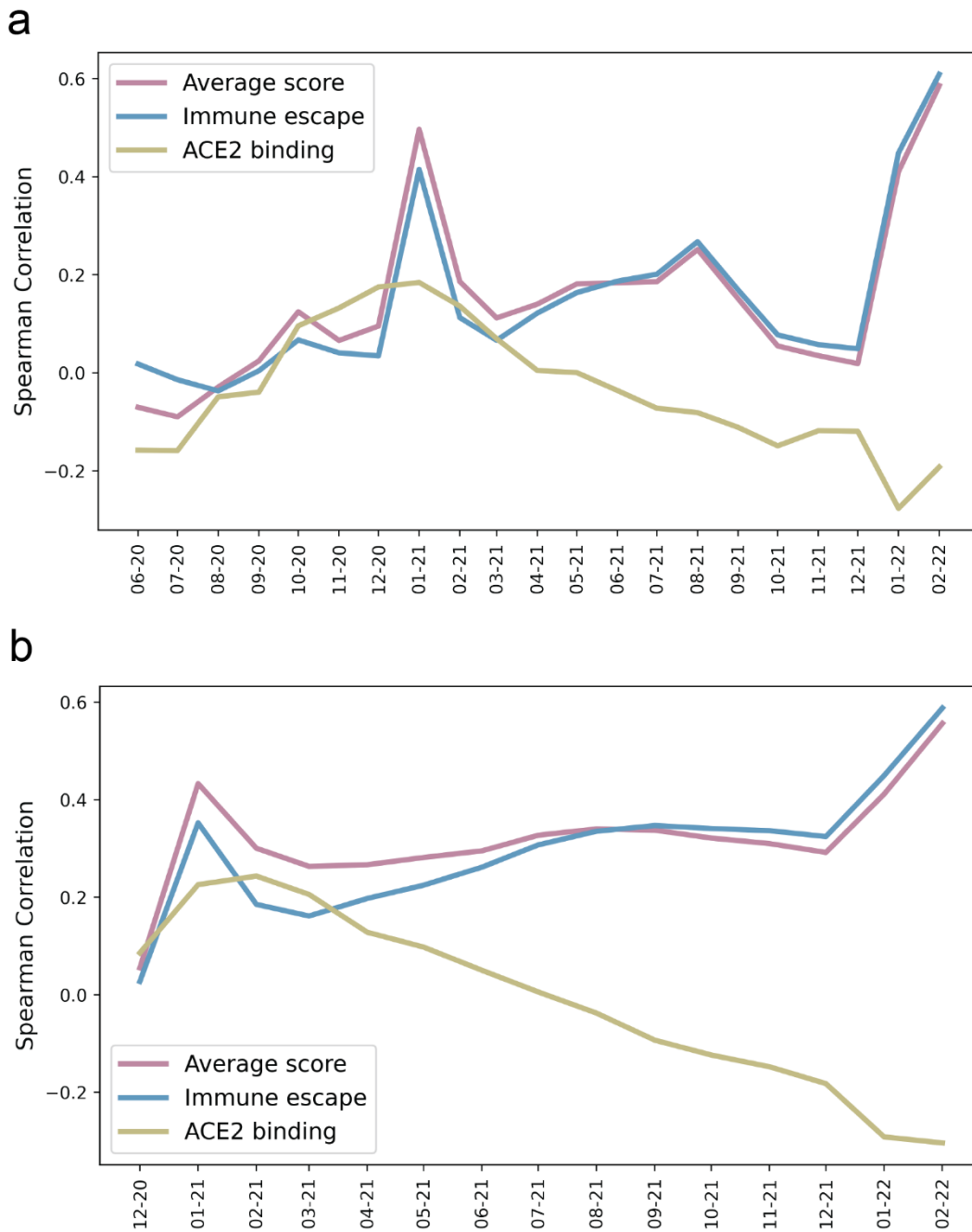


846

847 **Extended Data Fig. 3 | Pseudo time and ESM-1b model inference.** a, The landscape of RBD sequences,
848 represented as a KNN network and visualized use the UMAP, colored with the inferred pseudo time using our
849 model embeddings and scores. Gray lines indicate network edges. b and c, The landscape of the RBD sequences
850 from GISAID, represented as a KNN similarity network and visualized with the UMAP. The sequences
851 embeddings and directions among time are from the ESM-1b model. b, colored with the real-world sampling
852 time. Streamlines show the visual correlation between the ESM-1b inferred velocity and sampling time. c,
853 colored by the inferred pseudo time.

854

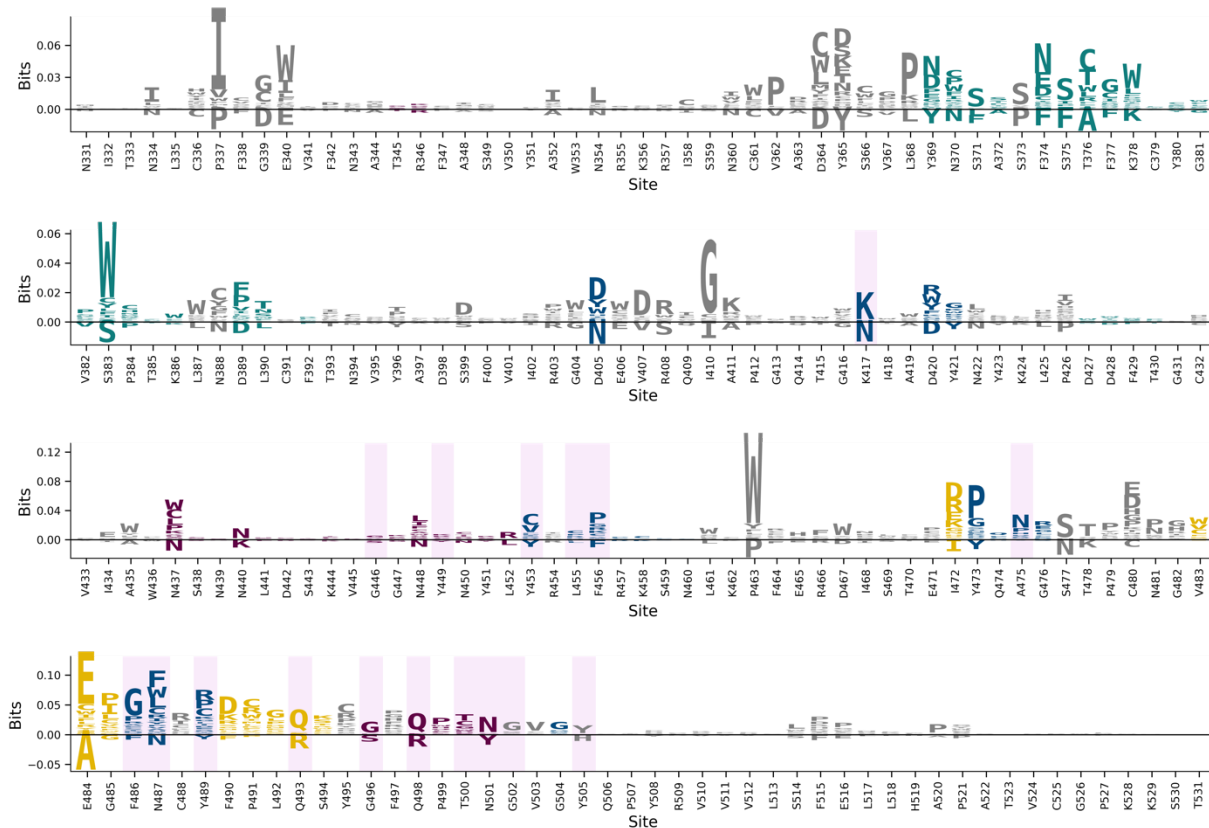
855



856

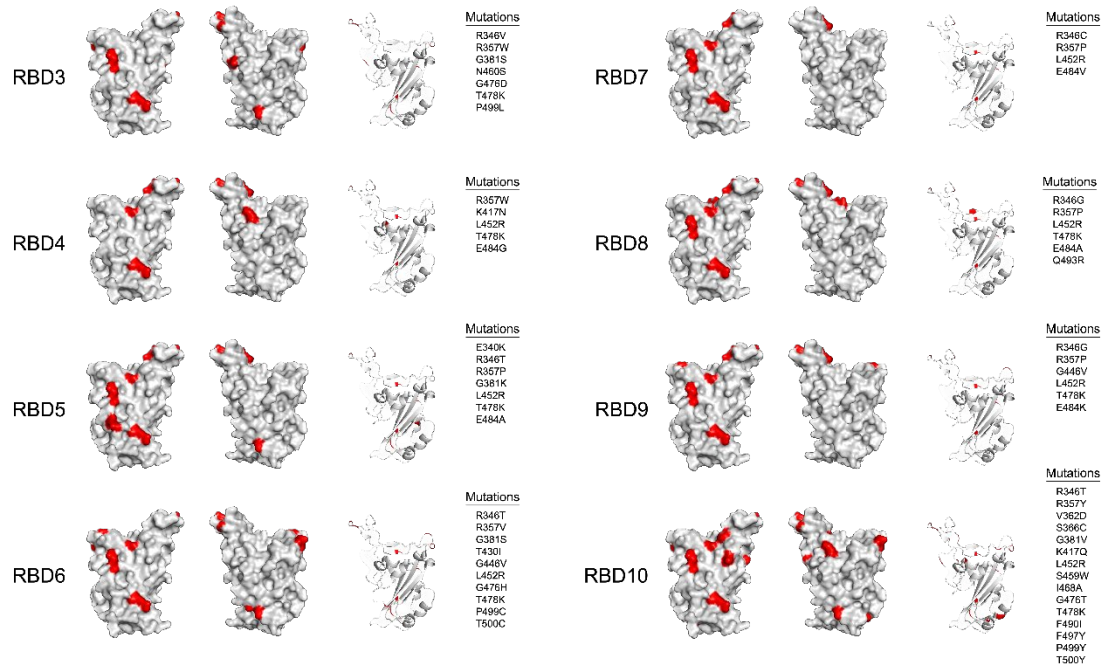
857 **Extended Data Fig. 4 | Correlation between the model scores and sampling date.** a, Spearman correlation
858 overtime for the model predictions, including the ACE2 binding score, immune escape potential, and the
859 weighted average of the two in a time window of previous six months for each sampled date. b, In a time window
860 of previous 12 months for each sampled date.

861



862
 863 **Extended Data Fig. 5 | KL logo plot for the entire RBD region.** Sequence logos generated from the
 864 differences of the generated sequences and the initial sequences, spanning the entire RBD region. The logos
 865 were calculated with probability weighted Kullback-Leibler divergence with a pseudo count of 0.1.

866
 867
 868
 869



870

871 **Extended Data Fig. 6 | Eight RBD mutants bearing different mutations on the surface were selected for**
872 **binding assay against mAbs.** Surface modeling of mutant RBD proteins was illustrated in grey. Mutations sites
873 were marked in red and listed beside the models.

874

875