

1 **COVID-19 coronavirus vaccine design using reverse vaccinology and machine**  
2 **learning**

3

4 **Edison Ong<sup>1</sup>, Mei U Wong<sup>2</sup>, Anthony Huffman<sup>1</sup>, Yongqun He<sup>1,2\*</sup>**

5

6 <sup>1</sup> Department of Computational Medicine and Bioinformatics, University of Michigan, Ann

7 Arbor, MI 48109, USA

8 <sup>2</sup> Unit for Laboratory Animal Medicine, Department of Microbiology and Immunology,

9 University of Michigan, Ann Arbor, MI 48109, USA

10

11

12 \*Corresponding authors:

13 Yongqun He: [yongqunh@med.umich.edu](mailto:yongqunh@med.umich.edu)

14

15

## 16 **Abstract**

17 To ultimately combat the emerging COVID-19 pandemic, it is desired to develop an  
18 effective and safe vaccine against this highly contagious disease caused by the SARS-CoV-2  
19 coronavirus. Our literature and clinical trial survey showed that the whole virus, as well as the  
20 spike (S) protein, nucleocapsid (N) protein, and membrane (M) protein, have been tested for  
21 vaccine development against SARS and MERS. However, these vaccine candidates might lack  
22 the induction of complete protection and have safety concerns. We then applied the Vaxign  
23 reverse vaccinology tool and the newly developed Vaxign-ML machine learning tool to predict  
24 COVID-19 vaccine candidates. By investigating the entire proteome of SARS-CoV-2, six  
25 proteins, including the S protein and five non-structural proteins (nsp3, 3CL-pro, and nsp8-10),  
26 were predicted to be adhesins, which are crucial to the viral adhering and host invasion. The S,  
27 nsp3, and nsp8 proteins were also predicted by Vaxign-ML to induce high protective  
28 antigenicity. Besides the commonly used S protein, the nsp3 protein has not been tested in any  
29 coronavirus vaccine studies and was selected for further investigation. The nsp3 was found to be  
30 more conserved among SARS-CoV-2, SARS-CoV, and MERS-CoV than among 15  
31 coronaviruses infecting human and other animals. The protein was also predicted to contain  
32 promiscuous MHC-I and MHC-II T-cell epitopes, and linear B-cell epitopes localized in specific  
33 locations and functional domains of the protein. By applying reverse vaccinology and machine  
34 learning, we predicted potential vaccine targets for effective and safe COVID-19 vaccine  
35 development. We then propose that an “Sp/Nsp cocktail vaccine” containing a structural  
36 protein(s) (Sp) and a non-structural protein(s) (Nsp) would stimulate effective complementary  
37 immune responses.

38

39

## 40 **Introduction**

41 The emerging Coronavirus Disease 2019 (COVID-19) pandemic poses a massive crisis to  
42 global public health. As of March 11, 2020, there were 118,326 confirmed cases and 4,292  
43 deaths, according to the World Health Organization (WHO), and WHO declared the COVID-19  
44 as a pandemic on the same day. As of March 22, there were >300,000 confirmed cases and  
45 >10,000 deaths globally in at least 167 countries, and the USA reported >27,000 confirmed cases

46 and >300 deaths. It is critical to develop an effective and safe vaccine(s) to control this fast-  
47 spreading disease and stop the pandemic.

48 The causative agent of the COVID-19 disease is the severe acute respiratory syndrome  
49 coronavirus 2 (SARS-CoV-2). Coronaviruses can cause animal diseases such as avian infectious  
50 bronchitis caused by the infectious bronchitis virus (IBV), and pig transmissible gastroenteritis  
51 caused by a porcine coronavirus<sup>1</sup>. Bats are commonly regarded as the natural reservoir of  
52 coronaviruses, which can be transmitted to humans and other animals after genetic mutations.  
53 There are seven known human coronaviruses, including the novel SARS-CoV-2. Four of them  
54 (HCoV-HKU1, HCoV-OC43, HCoV-229E, and HCoV-NL63) have been circulating in the  
55 human population worldwide and cause mild symptoms<sup>2</sup>. Coronavirus became prominence after  
56 Severe acute respiratory syndrome (SARS) and Middle East Respiratory Syndrome (MERS)  
57 outbreaks. In 2003, the SARS disease caused by the SARS-associated coronavirus (SARS-CoV)  
58 infected over 8,000 people worldwide and was contained in the summer of 2003<sup>3</sup>. SARS-CoV-2  
59 and SARS-CoV share high sequence identity<sup>4</sup>. The MERS disease infected more than 2,000  
60 people, which is caused by the MERS-associated coronavirus (MERS-CoV) and was first  
61 reported in Saudi Arabia and spread to several other countries since 2012<sup>5</sup>.

62 Although great efforts have been made to develop and manufacture COVID-19 vaccines,  
63 there is no human vaccine on the market to prevent this highly infectious disease. Coronaviruses  
64 are positively-stranded RNA viruses with its genome packed inside the nucleocapsid (N) protein  
65 and enveloped by the membrane (M) protein, envelope (E) protein, and the spike (S) protein<sup>6</sup>.  
66 While many coronavirus vaccine studies targeting different structural proteins were conducted,  
67 most of these efforts eventually ceased soon after the outbreak of SARS and MERS. With the  
68 recent COVID-19 pandemic outbreak, it is urgent to resume the coronavirus vaccine research. As  
69 the immediate response to the on-going pandemic, the first testing in humans of the mRNA-  
70 based vaccine targeting the S protein of SARS-CoV-2 (ClinicalTrials.gov Identifier:  
71 NCT04283461, Table 1) started on March 16, 2020. As the most superficial and protrusive  
72 protein of the coronaviruses, S protein plays a crucial role in mediating virus entry. In the SARS  
73 vaccine development, the full-length S protein and its S1 subunit (which contains receptor  
74 binding domain) have been frequently used as the vaccine antigens due to their ability to induce  
75 neutralizing antibodies that prevent host cell entry and infection.

76           However, the current coronavirus vaccines, including S protein-based vaccines, might  
77 have issues in the lack of inducing complete protection and possible safety concerns<sup>7,8</sup>. All  
78 existing SARS/MERS vaccines were reported to induce neutralizing antibodies and partial  
79 protection against the viral challenges in animal models (Table 2), but it is desired to induce  
80 complete protection or sterile immunity. Moreover, it has become increasingly clear that multiple  
81 immune responses, including those induced by humoral or cell-mediated immunity, are  
82 responsible for correlates of protection than antibody titers alone<sup>9</sup>. Both killed SARS-CoV whole  
83 virus vaccine and adenovirus-based recombinant vector vaccines expressing S or N proteins  
84 induced neutralizing antibody responses but did not provide complete protection in animal  
85 model<sup>10</sup>. A study has shown increased liver pathology in the vaccinated ferrets immunized with  
86 modified vaccinia Ankara-S recombinant vaccine<sup>11</sup>. The safety and efficacy of these vaccination  
87 strategies have not been fully tested in human clinical trials, but the safety can be a major  
88 concern. Therefore, novel strategies are needed to enhance the efficacy and safety of COVID-19  
89 vaccine development.

90           In recent years, the development of vaccine design has been revolutionized by the reverse  
91 vaccinology (RV), which aims to first identify promising vaccine candidate through  
92 bioinformatics analysis of the pathogen genome. RV has been successfully applied to vaccine  
93 discovery for pathogens such as Group B meningococcus and led to the license Bexsero  
94 vaccine<sup>12</sup>. Among current RV prediction tools<sup>13,14</sup>, Vaxign is the first web-based RV program<sup>15</sup>  
95 and has been used to successfully predict vaccine candidates against different bacterial and viral  
96 pathogens<sup>16-18</sup>. Recently we have also developed a machine learning approach called Vaxign-ML  
97 to enhance prediction accuracy<sup>19</sup>.

98           In this study, we first surveyed the existing coronavirus vaccine development status, and  
99 then applied the Vaxign RV and Vaxign-ML approaches to predict COVID-19 protein  
100 candidates for vaccine development. We identified six possible adhesins, including the structural  
101 S protein and five other non-structural proteins, and three of them (S, nsp3, and nsp8 proteins)  
102 were predicted to induce high protective immunity. The S protein was predicted to have the  
103 highest protective antigenicity score, and it has been extensively studied as the target of  
104 coronavirus vaccines by other researchers. The sequence conservation and immunogenicity of  
105 the multi-domain nsp3 protein, which was predicted to have the second-highest protective  
106 antigenicity score yet, was further analyzed in this study. Based on the predicted structural S

107 protein and non-structural proteins (including nsp3) using reverse vaccinology and machine  
108 learning, we proposed and discussed a cocktail vaccine strategy, for rational COVID-19 vaccine  
109 development.

110

## 111 **Results**

112

### 113 **Published research and clinical trial coronavirus vaccine studies**

114 To better understand the current status of coronavirus vaccine development, we  
115 systematically surveyed the development of vaccines for coronavirus from the ClinicalTrials.gov  
116 database and PubMed literature (as of March 17, 2020). Extensive effort has been made to  
117 develop a safe and effective vaccine against SARS or MERS, and the most advance clinical trial  
118 study is currently at phase II (Table 1). It is a challenging task to quickly develop a safe and  
119 effective vaccine for the on-going COVID-19 pandemic.

120 There are two primary design strategies for coronavirus vaccine development: the usage  
121 of the whole virus or genetically engineered vaccine antigens that can be delivered through  
122 different formats. The whole virus vaccines include inactivated<sup>20</sup> or live attenuated vaccines<sup>21,22</sup>  
123 (Table 2). The two live attenuated SARS vaccines mutated the exoribonuclease and envelop  
124 protein to reduce the virulence and/or replication capability of the SARS-CoV. Overall, the  
125 whole virus vaccines can induce a strong immune response and protect against coronavirus  
126 infections. Genetically engineered vaccines that target specific coronavirus protein are often used  
127 to improve vaccine safety and efficacy. The coronavirus antigens such as S protein, N protein,  
128 and M protein can be delivered as recombinant DNA vaccine and viral vector vaccine (Table 2).

129

### 130 **N protein is conserved among SARS-CoV-2, SARS-CoV, and MERS-CoV, but missing from** 131 **the other four human coronaviruses causing mild symptoms**

132 We first used the Vaxign analysis framework<sup>15,19</sup> to compare the full proteomes of seven  
133 human coronavirus strains (SARS-CoV-2, SARS-CoV, MERS-CoV, HCoV-229E, HCoV-  
134 OC43, HCoV-NL63, and HCoV-HKU1). The proteins of SARS-CoV-2 were used as the seed for  
135 the pan-genomic comparative analysis. The Vaxign pan-genomic analysis reported only the N  
136 protein in SARS-CoV-2 having high sequence similarity among the more severe form of  
137 coronavirus (SARS-CoV and MERS-CoV), while having low sequence similarity among the

138 more typically mild HCoV-229E, HCoV-OC43, HCoV-NL63, and HCoV-HKU1. The sequence  
139 conservation suggested the potential of N protein as a candidate for the cross-protective vaccine  
140 against SARS and MERS. The N protein was also evaluated and used for vaccine development  
141 (Table 2). The N protein packs the coronavirus RNA to form the helical nucleocapsid in virion  
142 assembly. This protein is more conserved than the S protein and was reported to induce an  
143 immune response and neutralize coronavirus infections<sup>23</sup>. However, a study also showed the  
144 linkage between N protein and severe pneumonia or other serious liver failures related to the  
145 pathogenesis of SARS<sup>24</sup>.

146

### 147 **Six adhesive proteins in SARS-CoV-2 identified as potential vaccine targets**

148 The Vaxign RV analysis predicted six SARS-CoV-2 proteins (S protein, nsp3, 3CL-PRO,  
149 and nsp8-10) as adhesive proteins (Table 3). Adhesin plays a critical role in the virus adhering to  
150 the host cell and facilitating the virus entry to the host cell<sup>25</sup>, which has a significant association  
151 with the vaccine-induced protection<sup>26</sup>. In SARS-CoV-2, S protein was predicted to be adhesin,  
152 matching its primary role in virus entry. The structure of SARS-CoV-2 S protein was determined<sup>27</sup>  
153 and reported to contribute to the host cell entry by interacting with the angiotensin-converting  
154 enzyme 2 (ACE2)<sup>28</sup>. Besides S protein, the other five predicted adhesive proteins were all non-  
155 structural proteins. In particular, nsp3 is the largest non-structural protein of SARS-CoV-2  
156 comprises various functional domains<sup>29</sup>.

157

### 158 **Three adhesin proteins were predicted to induce strong protective immunity**

159 The Vaxign-ML pipeline computed the protegenicity (protective antigenicity) score and  
160 predicted the induction of protective immunity by a vaccine candidate<sup>19</sup>. The training data  
161 consisted of viral protective antigens, which were tested to be protective in at least one animal  
162 challenge model<sup>30</sup>. The performance of the Vaxign-ML models was evaluated (Table S1 and  
163 Figure S1), and the best performing model had a weighted F1-score of 0.94. Using the optimized  
164 Vaxign-ML model, we predicted three proteins (S protein, nsp3, and nsp8) as vaccine candidates  
165 with significant protegenicity scores (Table 3). The S protein was predicted to have the highest  
166 protegenicity score, which is consistent with the experimental observations reported in the  
167 literature. The nsp3 protein is the second most promising vaccine candidate besides S protein.  
168 There was currently no study of nsp3 as a vaccine target. The structure and functions of this protein

169 have various roles in coronavirus infection, including replication and pathogenesis (immune  
170 evasion and virus survival)<sup>29</sup>. Therefore, we selected nsp3 for further investigation, as described  
171 below.

172

### 173 **Nsp3 as a vaccine candidate**

174 The multiple sequence alignment and the resulting phylogeny of nsp3 protein showed that  
175 this protein in SARS-CoV-2 was more closely related to the human coronaviruses SARS-CoV and  
176 MERS-CoV, and bat coronaviruses BtCoV/HKU3, BtCoV/HKU4, and BtCoV/HKU9. We studied  
177 the genetic conservation of nsp3 protein (Figure 1A) in seven human coronaviruses and eight  
178 coronaviruses infecting other animals (Table S2). The five human coronaviruses, SARS-CoV-2,  
179 SARS-CoV, MERS-CoV, HCoV-HKU1, and HCoV-OC43, belong to the beta-coronavirus while  
180 HCoV-229E and HCoV-NL63 belong to the alpha-coronavirus. The HCoV-HKU1 and HCoV-  
181 OC43, as the human coronavirus with mild symptoms clustered together with murine MHV-A59.  
182 The more severe form of human coronavirus SARS-CoV-2, SARS-CoV, and MERS-CoV grouped  
183 with three bat coronaviruses BtCoV/HKU3, BtCoV/HKU4, and BtCoV/HKU9.

184 When evaluating the amino acid conservations relative to the functional domains in nsp3,  
185 all protein domains, except the hypervariable region (HVR), macro-domain 1 (MAC1) and beta-  
186 coronavirus-specific marker  $\beta$ SM, showed higher conservation in SARS-CoV-2, SARS-CoV, and  
187 MERS-CoV (Figure 1B). The amino acid conservation between the major human coronavirus  
188 (SARS-CoV-2, SARS-CoV, and MERS-CoV) was plotted and compared to all 15 coronaviruses  
189 used to generate the phylogenetic of nsp3 protein (Figure 1B). The SARS-CoV domains were also  
190 plotted (Figure 1B), with the relative position in the multiple sequence alignment (MSA) of all 15  
191 coronaviruses (Table S3 and Figure S2).

192 The immunogenicity of nsp3 protein in terms of T cell MHC-I & MHC-II and linear B cell  
193 epitopes was also investigated. There were 28 and 42 promiscuous epitopes predicted to bind the  
194 reference MHC-I & MHC-II alleles, which covered the majority of the world population,  
195 respectively (Table S4-5). In terms of linear B cell epitopes, there were 14 epitopes with BepiPred  
196 scores over 0.55 and had at least ten amino acids in length (Table S6). The 3D structure of SARS-  
197 CoV-2 protein was plotted and highlighted with the T cell MHC-I & MHC-II, and linear B cell  
198 epitopes (Figure 2). The predicted B cell epitopes were more likely located in the distal region of  
199 the nsp3 protein structure. Most of the predicted MHC-I & MHC-II epitopes were embedded inside

200 the protein. The sliding averages of T cell MHC-I & MHC-II and linear B cell epitopes were  
201 plotted with respect to the tentative SARS-CoV-2 nsp3 protein domains using SARS-CoV nsp3  
202 protein as a reference (Figure 3). The ubiquitin-like domain 1 and 2 (Ubl1 and Ubl2) only predicted  
203 to have MHC-I epitopes. The Domain Preceding Ubl2 and PL2-PRO (DPUP) domain had only  
204 predicted MHC-II epitopes. The PL2-PRO contained both predicted MHC-I and MHC-II epitopes,  
205 but not B cell epitopes. In particular, the TM1, TM2, and AH1 were predicted helical regions with  
206 high T cell MHC-I and MHC-II epitopes<sup>31</sup>. The TM1 and TM2 are transmembrane regions passing  
207 the endoplasmic reticulum (ER) membrane. The HVR, MAC2, MAC3, nucleic-acid binding  
208 domain (NAB),  $\beta$ SM, Nsp3 ectodomain; (3Ecto), Y1, and CoV-Y domain contained predicted B  
209 cell epitopes. Finally, the Vaxign RV framework also predicted 2 regions (position 251-260 and  
210 329-337) in the MAC1 domain of nsp3 domain having high sequence similarity to the human  
211 mono-ADP-ribosyltransferase PARP14 (NP\_060024.2).

212

## 213 Discussion

214 Our prediction of the potential SARS-CoV-2 antigens, which could induce protective  
215 immunity, provides a timely analysis for the vaccine development against COVID-19. Currently,  
216 most coronavirus vaccine studies use the whole inactivated or attenuated virus, or target the  
217 structural proteins such as the spike (S) protein, nucleocapsid (N) protein, and membrane (M)  
218 protein (Table 2). But the inactivated or attenuated whole virus vaccine might induce strong  
219 adverse events. On the other hand, vaccines targeting the structural proteins induce a strong  
220 immune response<sup>23,32,33</sup>. In some studies, these structural proteins, including the S and N proteins,  
221 were reported to associate with the pathogenesis of coronavirus<sup>24,34</sup> and might raise safety  
222 concern<sup>11</sup>. Our study applied state-of-the-art Vaxign reserve vaccinology (RV) and Vaxign-ML  
223 machine learning strategies to the entire SARS-CoV-2 proteomes, including both structural and  
224 non-structural proteins for vaccine candidate prediction. Our results indicate, for the first time, that  
225 many non-structural proteins could be used as potential vaccine candidates.

226 The SARS-CoV-2 S protein was identified by our Vaxign and Vaxign-ML analysis as the  
227 most favorable vaccine candidate. First, the Vaxign RV framework predicted the S protein as a  
228 likely adhesin, which is consistent with the role of S protein for the invasion of host cells. Second,  
229 our Vaxign-ML predicted that the S protein had a high protective antigenicity score. These results  
230 confirmed the role of S protein as the important target of COVID-19 vaccines. However, targeting



231 only the S protein may induce high serum-neutralizing antibody titers but cannot induce complete  
232 protection<sup>10</sup>. In addition, HCoV-NL63 also uses S protein and employs the angiotensin-converting  
233 enzyme 2 (ACE2) for cellular entry, despite markedly weak pathogenicity<sup>35</sup>. This suggests that the  
234 S protein is not the only factor determining the infection level of a human coronavirus. Thus,  
235 alternative vaccine antigens may be considered as potential targets for COVID-19 vaccines.

236 Among the five non-structural proteins being predicted as potential vaccine candidates, the  
237 nsp3 protein was predicted to have second-highest protective antigenicity score, adhesin property,  
238 promiscuous MHC-I & MHC-II T cell epitopes, and B cell epitopes. The nsp3 is the largest non-  
239 structural protein that includes multiple functional domains related to viral pathogenesis<sup>29</sup>. The  
240 multiple sequence alignment of nsp3 also showed higher sequence conservation in most of the  
241 functional domains in SARS-CoV-2, SARS-CoV, and MERS-CoV, than in all 15 coronavirus  
242 strains (Fig. 1B). Besides the nsp3 protein, our study also predicted four additional non-structural  
243 proteins (3CL-pro, nsp8, nsp9, and nsp10) as possible vaccine candidates based on their adhesin  
244 probabilities, and the nsp8 protein was also predicted to have a significant protective antigenicity  
245 score.

246 However, these predicted non-structural proteins (nsp3, 3CL-pro, nsp8, nsp9, and nsp10)  
247 are not part of the viral structural particle, and all the current SARS/MERS/COVID-19 vaccine  
248 studies target the structural (S/M/N) proteins. Although structural proteins are commonly used as  
249 viral vaccine candidates, non-structural proteins correlates to vaccine protection. The non-  
250 structural protein NS1 was found to induce protective immunity against the infections by  
251 flaviviruses<sup>36</sup>. Since NS1 is not part of the virion, antibodies against NS1 have no neutralizing  
252 activity but some exhibit complement-fixing activity<sup>37</sup>. However, passive transfer of anti-NS1  
253 antibody or immunization with NS1 conferred protection<sup>38</sup>. Anti-NS1 antibody could also reduce  
254 viral replication by complement-dependent cytotoxicity of infected cells, block NS1-induced  
255 pathogenic effects, and attenuate NS1-induced disease development during the critical phase<sup>39</sup>.  
256 Finally, NS1 is not a structural protein and anti-NS1 antibody will not induce antibody-dependent  
257 enhancement (ADE), which is a virulence factor and a risk factor causing many adverse events<sup>39</sup>.  
258 The non-structural proteins of the hepatitis C virus were reported to induce HCV-specific vigorous  
259 and broad-spectrum T-cell responses<sup>40</sup>. The non-structural HIV-1 gene products were also shown  
260 to be valuable targets for prophylactic or therapeutic vaccines<sup>41</sup>. Therefore, it is reasonable to  
261 consider the SARS-CoV-2 non-structural proteins (e.g., nsp3) as possible vaccine targets, which

262 might induce cell-mediated or humoral immunity necessary to prevent viral invasion and/or  
263 replication. None of the non-structural proteins have been evaluated as vaccine candidates, and the  
264 feasibility of these proteins as vaccine targets are subject to further experimental verification.

265 In addition to vaccines expressing a single or a combination of structural proteins, here we  
266 propose an “Sp/Nsp cocktail vaccine” as an effective strategy for COVID-19 vaccine development.  
267 A typical cocktail vaccine includes more than one antigen to cover different aspects of  
268 protection<sup>42,43</sup>. The licensed Group B meningococcus Bexsero vaccine, which was developed via  
269 reverse vaccinology, contains three protein antigens<sup>12</sup>. To develop an efficient and safe COVID-  
270 19 cocktail vaccine, an “Sp/Nsp cocktail vaccine”, which mixes a structural protein(s) (Sp, such  
271 as S protein) and a non-structural protein(s) (Nsp, such as nsp3) could induce more favorable  
272 protective immune responses than vaccines expressing a structural protein(s). The benefit of a  
273 cocktail vaccine strategy could induce immunity that can protect the host against not only the S-  
274 ACE2 interaction and viral entry to the host cells, but also protect against the accessory non-  
275 structural adhesin proteins (e.g., nsp3), which might also be vital to the viral entry and replication.  
276 The usage of more than one antigen allows us to reduce the volume of each antigen and thus to  
277 reduce the induction of adverse events. Nonetheless, the potentials of the proposed “Sp/Nsp  
278 cocktail vaccine” strategy need to be experimentally validated.

279 For rational COVID-19 vaccine development, it is critical to understand the fundamental  
280 host-coronavirus interaction and protective immune mechanism<sup>7</sup>. Such understanding may not  
281 only provide us guidance in terms of antigen selection but also facilitate our design of vaccine  
282 formulations. For example, an important foundation of our prediction in this study is based on our  
283 understanding of the critical role of adhesin as a virulence factor as well as protective antigen. The  
284 choice of DNA vaccine, recombinant vaccine vector, and another method of vaccine formulation  
285 is also deeply rooted in our understanding of pathogen-specific immune response induction.  
286 Different experimental conditions may also affect results<sup>44,45</sup>. Therefore, it is crucial to understand  
287 the underlying molecular and cellular mechanisms for rational vaccine development.

288

## 289 **Methods**

290 **Annotation of literature and database records.** We annotated peer-reviewed journal articles  
291 stored in the PubMed database and the ClinicalTrials.gov database. From the peer-reviewed  
292 articles, we identified and annotated those coronavirus vaccine candidates that were

293 experimentally studied and found to induce protective neutralizing antibody or provided immunity  
294 against virulent pathogen challenge.

295

296 **Vaxign prediction.** The SARS-CoV-2 sequence was obtained from NCBI. All the proteins of six  
297 known human coronavirus strains, including SARS-CoV, MERS-CoV, HCoV-229E, HCoV-  
298 OC43, HCoV-NL63, and HCoV-HKU1 were extracted from Uniprot proteomes<sup>46</sup>. The full  
299 proteomes of these seven coronaviruses were then analyzed using the Vaxign reverse vaccinology  
300 pipeline<sup>15,19</sup>. The Vaxign program predicted several biological features, including adhesin  
301 probability<sup>47</sup>, transmembrane helix<sup>48</sup>, orthologous proteins<sup>49</sup>, and protein functions<sup>15,19</sup>.

302

303 **Vaxign-ML prediction.** The ML-based RV prediction model was built following a similar  
304 methodology described in the Vaxign-ML<sup>19</sup>. Specifically, the positive samples in the training data  
305 included 397 bacterial and 178 viral protective antigens (PAgs) recorded in the Protegen database<sup>30</sup>  
306 after removing homologous proteins with over 30% sequence identity. There were 4,979 negative  
307 samples extracted from the corresponding pathogens' Uniprot proteomes<sup>46</sup> with sequence dis-  
308 similarity to the PAgs, as described in previous studies<sup>50-52</sup>. Homologous proteins in the negative  
309 samples were also removed. The proteins in the resulting dataset were annotated with biological  
310 and physicochemical features. The biological features included adhesin probability<sup>47</sup>,  
311 transmembrane helix<sup>48</sup>, and immunogenicity<sup>53</sup>. The physicochemical features included the  
312 compositions, transitions and distributions<sup>54</sup>, quasi-sequence-order<sup>55</sup>, Moreau-Broto auto-  
313 correlation<sup>56,57</sup>, and Geary auto-correlation<sup>58</sup> of various physicochemical properties such as charge,  
314 hydrophobicity, polarity, and solvent accessibility<sup>59</sup>. Five supervised ML classification algorithms,  
315 including logistic regression, support vector machine, k-nearest neighbor, random forest<sup>60</sup>, and  
316 extreme gradient boosting (XGB)<sup>61</sup> were trained on the annotated proteins dataset. The  
317 performance of these models was evaluated using a nested five-fold cross-validation (N5CV)  
318 based on the area under receiver operating characteristic curve, precision, recall, weighted F1-  
319 score, and Matthew's correlation coefficient. The best performing XGB model was selected to  
320 predict the protegenicity score of all SARS-CoV-2 isolate Wuhan-Hu-1 (GenBank ID:  
321 MN908947.3) proteins, downloaded from NCBI. A protein with protegenicity score over 0.9 is  
322 considered as strong vaccine candidate (weighted F1-score > 0.94 in N5CV).

323

324 **Phylogenetic analysis.** The protein nsp3 was selected for further investigation. The nsp3 proteins  
325 of 14 coronaviruses besides SARS-CoV-2 were downloaded from the Uniprot (Table S2). Multiple  
326 sequence alignment of these nsp3 proteins was performed using MUSCLE<sup>62</sup> and visualized via  
327 SEAVIEW<sup>63</sup>. The phylogenetic tree was constructed using PhyML<sup>64</sup>, and the amino acid  
328 conservation was estimated by the Jensen-Shannon Divergence (JSD)<sup>65</sup>. The JSD score was also  
329 used to generate a sequence conservation line using the nsp3 protein sequences from 4 or 13  
330 coronaviruses.

331  
332 **Immunogenicity analysis.** The immunogenicity of the nsp3 protein was evaluated by the  
333 prediction of T cell MHC-I and MHC-II, and linear B cell epitopes. For T cell MHC-I epitopes,  
334 the IEDB consensus method was used to predicting promiscuous epitopes binding to 4 out of 27  
335 MHC-I reference alleles with consensus percentile ranking less than 1.0 score<sup>53</sup>. For T cell MHC-  
336 II epitopes, the IEDB consensus method was used to predicting promiscuous epitopes binding to  
337 more than half of the 27 MHC-II reference alleles with consensus percentile ranking less than 10.0.  
338 The MHC-I and MHC-II reference alleles covered a wide range of human genetic variation  
339 representing the majority of the world population<sup>66,67</sup>. The linear B cell epitopes were predicted  
340 using the BepiPred 2.0 with a cutoff of 0.55 score<sup>68</sup>. Linear B cell epitopes with at least ten amino  
341 acids were mapped to the predicted 3D structure of SARS-CoV-2 nsp3 protein visualized via  
342 PyMol<sup>69</sup>. The predicted count of T cell MHC-I and MHC-II epitopes, and the predicted score of  
343 linear B cell epitopes were computed as the sliding averages with a window size of ten amino acids.  
344 The nsp3 protein 3D structure was predicted using C-I-Tasser<sup>70</sup> available in the Zhang Lab  
345 webserver (<https://zhanglab.ccmb.med.umich.edu/C-I-TASSER/2019-nCov/>).

346

## 347 **References**

- 348 1. Perlman, S. & Netland, J. Coronaviruses post-SARS: Update on replication and  
349 pathogenesis. *Nature Reviews Microbiology* (2009). doi:10.1038/nrmicro2147
- 350 2. Cabeça, T. K., Granato, C. & Bellei, N. Epidemiological and clinical features of human  
351 coronavirus infections among different subsets of patients. *Influenza Other Respi. Viruses*  
352 (2013). doi:10.1111/irv.12101
- 353 3. Lu, R. *et al.* Genomic characterisation and epidemiology of 2019 novel coronavirus:  
354 implications for virus origins and receptor binding. *Lancet* (2020). doi:10.1016/S0140-

- 355 6736(20)30251-8
- 356 4. Lai, C.-C., Shih, T.-P., Ko, W.-C., Tang, H.-J. &Hsueh, P.-R. Severe acute respiratory  
357 syndrome coronavirus 2 (SARS-CoV-2) and coronavirus disease-2019 (COVID-19): The  
358 epidemic and the challenges. *Int. J. Antimicrob. Agents* (2020).  
359 doi:10.1016/j.ijantimicag.2020.105924
- 360 5. Chan, J. F. W. *et al.* Middle East Respiratory syndrome coronavirus: Another zoonotic  
361 betacoronavirus causing SARS-like disease. *Clin. Microbiol. Rev.* (2015).  
362 doi:10.1128/CMR.00102-14
- 363 6. Li, F. Structure, Function, and Evolution of Coronavirus Spike Proteins. *Annu. Rev. Virol.*  
364 (2016). doi:10.1146/annurev-virology-110615-042301
- 365 7. Roper, R. L. &Rehm, K. E. SARS vaccines: Where are we? *Expert Review of Vaccines*  
366 (2009). doi:10.1586/erv.09.43
- 367 8. deWit, E., vanDoremalen, N., Falzarano, D. &Munster, V. J. SARS and MERS: recent  
368 insights into emerging coronaviruses. *Nat. Rev. Microbiol.* **14**, 523–534 (2016).
- 369 9. Plotkin, S. A. Updates on immunologic correlates of vaccine-induced protection. *Vaccine*  
370 **38**, 2250–2257 (2020).
- 371 10. See, R. H. *et al.* Severe acute respiratory syndrome vaccine efficacy in ferrets: Whole  
372 killed virus and adenovirus-vectored vaccines. *J. Gen. Virol.* (2008).  
373 doi:10.1099/vir.0.2008/001891-0
- 374 11. Weingartl, H. *et al.* Immunization with Modified Vaccinia Virus Ankara-Based  
375 Recombinant Vaccine against Severe Acute Respiratory Syndrome Is Associated with  
376 Enhanced Hepatitis in Ferrets. *J. Virol.* (2004). doi:10.1128/jvi.78.22.12672-12676.2004
- 377 12. Folaranmi, T., Rubin, L., Martin, S. W., Patel, M. &MacNeil, J. R. Use of Serogroup B  
378 Meningococcal Vaccines in Persons Aged  $\geq 10$  Years at Increased Risk for Serogroup B  
379 Meningococcal Disease: Recommendations of the Advisory Committee on Immunization  
380 Practices, 2015. *MMWR Morb Mortal Wkly Rep* **64**, 608–612 (2015).
- 381 13. He, Y. *et al.* Emerging vaccine informatics. *J. Biomed. Biotechnol.* **2010**, (2010).
- 382 14. Dalsass, M., Brozzi, A., Medini, D. &Rappuoli, R. Comparison of Open-Source Reverse  
383 Vaccinology Programs for Bacterial Vaccine Antigen Discovery. *Front. Immunol.* **10**, 1–  
384 12 (2019).
- 385 15. He, Y., Xiang, Z. &Mobley, H. L. T. Vaxign: The first web-based vaccine design program

- 386 for reverse vaccinology and applications for vaccine development. *J. Biomed. Biotechnol.*  
387 **2010**, (2010).
- 388 16. Xiang, Z. A. &He, Y. O. Genome-wide prediction of vaccine targets for human herpes  
389 simplex viruses using Vaxign reverse vaccinology Human Herpes Simplex ( HSV )  
390 Viruses. **14**, 1–10 (2013).
- 391 17. Singh, R., Garg, N., Shukla, G., Capalash, N. &Sharma, P. Immunoprotective Efficacy of  
392 Acinetobacter baumannii Outer Membrane Protein, FilF, Predicted In silico as a Potential  
393 Vaccine Candidate. *Front. Microbiol.* **7**, (2016).
- 394 18. Navarro-Quiroz, E. *et al.* Prediction of Epitopes in the Proteome of Helicobacter pylori.  
395 *Glob. J. Health Sci.* **10**, 148 (2018).
- 396 19. Ong, E. *et al.* Vaxign-ML: Supervised Machine Learning Reverse Vaccinology Model for  
397 Improved Prediction of Bacterial Protective Antigens. *Bioinformatics* (2020).
- 398 20. See, R. H. *et al.* Comparative evaluation of two severe acute respiratory syndrome  
399 (SARS) vaccine candidates in mice challenged with SARS coronavirus. *J. Gen. Virol.*  
400 (2006). doi:10.1099/vir.0.81579-0
- 401 21. Graham, R. L. *et al.* A live, impaired-fidelity coronavirus vaccine protects in an aged,  
402 immunocompromised mouse model of lethal disease. *Nat. Med.* (2012).  
403 doi:10.1038/nm.2972
- 404 22. Fett, C., DeDiego, M. L., Regla-Nava, J. A., Enjuanes, L. &Perlman, S. Complete  
405 Protection against Severe Acute Respiratory Syndrome Coronavirus-Mediated Lethal  
406 Respiratory Disease in Aged Mice by Immunization with a Mouse-Adapted Virus Lacking  
407 E Protein. *J. Virol.* (2013). doi:10.1128/jvi.00087-13
- 408 23. Zhao, P. *et al.* Immune responses against SARS-coronavirus nucleocapsid protein induced  
409 by DNA vaccine. *Virology* (2005). doi:10.1016/j.virol.2004.10.016
- 410 24. Yasui, F. *et al.* Prior Immunization with Severe Acute Respiratory Syndrome (SARS)-  
411 Associated Coronavirus (SARS-CoV) Nucleocapsid Protein Causes Severe Pneumonia in  
412 Mice Infected with SARS-CoV. *J. Immunol.* (2008). doi:10.4049/jimmunol.181.9.6337
- 413 25. Ribet, D. &Cossart, P. How bacterial pathogens colonize their hosts and invade deeper  
414 tissues. *Microbes Infect.* **17**, 173–183 (2015).
- 415 26. Ong, E., Wong, M. U. &He, Y. Identification of New Features from Known Bacterial  
416 Protective Vaccine Antigens Enhances Rational Vaccine Design. *Front. Immunol.* **8**, 1–11

- 417 (2017).
- 418 27. Wrapp, D. *et al.* Cryo-EM structure of the 2019-nCoV spike in the prefusion  
419 conformation. *Science* (2020). doi:10.1126/science.abb2507
- 420 28. Letko, M., Marzi, A. & Munster, V. Functional assessment of cell entry and receptor usage  
421 for SARS-CoV-2 and other lineage B betacoronaviruses. *Nat. Microbiol.* (2020).  
422 doi:10.1038/s41564-020-0688-y
- 423 29. Lei, J., Kusov, Y. & Hilgenfeld, R. Nsp3 of coronaviruses: Structures and functions of a  
424 large multi-domain protein. *Antiviral Research* **149**, 58–74 (2018).
- 425 30. Yang, B., Sayers, S., Xiang, Z. & He, Y. Protegen: A web-based protective antigen  
426 database and analysis system. *Nucleic Acids Res.* **39**, 1073–1078 (2011).
- 427 31. Rothbard, J. B. & Taylor, W. R. A sequence pattern common to T cell epitopes. *EMBO J.*  
428 (1988). doi:10.1002/j.1460-2075.1988.tb02787.x
- 429 32. Shi, S. Q. *et al.* The expression of membrane protein augments the specific responses  
430 induced by SARS-CoV nucleocapsid DNA immunization. *Mol. Immunol.* (2006).  
431 doi:10.1016/j.molimm.2005.11.005
- 432 33. Al-Amri, S. S. *et al.* Immunogenicity of Candidate MERS-CoV DNA Vaccines Based on  
433 the Spike Protein. *Sci. Rep.* (2017). doi:10.1038/srep44875
- 434 34. Glansbeek, H. L. *et al.* Adverse effects of feline IL-12 during DNA vaccination against  
435 feline infectious peritonitis virus. *J. Gen. Virol.* (2002). doi:10.1099/0022-1317-83-1-1
- 436 35. Hofmann, H. *et al.* Human coronavirus NL63 employs the severe acute respiratory  
437 syndrome coronavirus receptor for cellular entry. *Proc. Natl. Acad. Sci. U. S. A.* (2005).  
438 doi:10.1073/pnas.0409465102
- 439 36. Salat, J. *et al.* Tick-borne encephalitis virus vaccines contain non-structural protein 1  
440 antigen and may elicit NS1-specific antibody responses in vaccinated individuals.  
441 *Vaccines* (2020). doi:10.3390/vaccines8010081
- 442 37. Schlesinger, J. J., Brandriss, M. W. & Walsh, E. E. Protection against 17D yellow fever  
443 encephalitis in mice by passive transfer of monoclonal antibodies to the nonstructural  
444 glycoprotein gp48 and by active immunization with gp48. *J. Immunol.* (1985).
- 445 38. Gibson, C. A., Schlesinger, J. J. & Barrett, A. D. T. Prospects for a virus non-structural  
446 protein as a subunit vaccine. *Vaccine* (1988). doi:10.1016/0264-410X(88)90004-7
- 447 39. Chen, H. R., Lai, Y. C. & Yeh, T. M. Dengue virus non-structural protein 1: A pathogenic

- 448 factor, therapeutic target, and vaccine candidate. *Journal of Biomedical Science* (2018).  
449 doi:10.1186/s12929-018-0462-0
- 450 40. Ip, P. P. *et al.* Alphavirus-based vaccines encoding nonstructural proteins of hepatitis c  
451 virus induce robust and protective T-cell responses. *Mol. Ther.* (2014).  
452 doi:10.1038/mt.2013.287
- 453 41. Cafaro, A. *et al.* Anti-tat immunity in HIV-1 infection: Effects of naturally occurring and  
454 vaccine-induced antibodies against tat on the course of the disease. *Vaccines* (2019).  
455 doi:10.3390/vaccines7030099
- 456 42. Sealy, R. *et al.* Preclinical and clinical development of a multi-envelope, DNA-virus-  
457 protein (D-V-P) HIV-1 vaccine. *International Reviews of Immunology* (2009).  
458 doi:10.1080/08830180802495605
- 459 43. Millet, P. *et al.* Immunogenicity of the Plasmodium falciparum asexual blood-stage  
460 synthetic peptide vaccine SPf66. *Am. J. Trop. Med. Hyg.* (1993).  
461 doi:10.4269/ajtmh.1993.48.424
- 462 44. He, Y. *et al.* Updates on the web-based VIOLIN vaccine database and analysis system.  
463 *Nucleic Acids Res.* **42**, 1124–1132 (2014).
- 464 45. Ong, E. *et al.* VIO: Ontology classification and study of vaccine responses given various  
465 experimental and analytical conditions. *BMC Bioinformatics* (2019). doi:10.1186/s12859-  
466 019-3194-6
- 467 46. The UniProt Consortium. The Universal Protein Resource (UniProt). *Nucleic Acids Res.*  
468 **36**, D193-7 (2008).
- 469 47. Sachdeva, G., Kumar, K., Jain, P. & Ramachandran, S. SPAAN: A software program for  
470 prediction of adhesins and adhesin-like proteins using neural networks. *Bioinformatics* **21**,  
471 483–491 (2005).
- 472 48. Krogh, A., Larsson, B., vonHeijne, G. & Sonnhammer, E. L. . Predicting transmembrane  
473 protein topology with a hidden Markov model: application to complete genomes. *J Mol*  
474 *Biol* **305**, 567–580 (2001).
- 475 49. Li, L., Stoeckert, C. J. & Roos, D. S. OrthoMCL: Identification of ortholog groups for  
476 eukaryotic genomes. *Genome Res.* (2003). doi:10.1101/gr.1224503
- 477 50. Bowman, B. N. *et al.* Improving reverse vaccinology with a machine learning approach.  
478 *Vaccine* **29**, 8156–8164 (2011).



- 479 51. Doytchinova, I. a & Flower, D. R. VaxiJen: a server for prediction of protective antigens,  
480 tumour antigens and subunit vaccines. *BMC Bioinformatics* **8**, 4 (2007).
- 481 52. Heinson, A. I. *et al.* Enhancing the biological relevance of machine learning classifiers for  
482 reverse vaccinology. *Int. J. Mol. Sci.* **18**, (2017).
- 483 53. Fleri, W. *et al.* The immune epitope database and analysis resource in epitope discovery  
484 and synthetic vaccine design. *Front. Immunol.* **8**, 1–16 (2017).
- 485 54. Dubchak, I., Muchnik, I., Holbrook, S. R. & Kim, S. H. Prediction of protein folding class  
486 using global description of amino acid sequence. *Proc. Natl. Acad. Sci. U. S. A.* **92**, 8700–  
487 8704 (1995).
- 488 55. Chou, K.-C. Prediction of Protein Subcellular Locations by Incorporating Quasi-  
489 Sequence-Order Effect. *Biochem. Biophys. Res. Commun.* **278**, 477–483 (2000).
- 490 56. Lin, Z. & Pan, X. M. Accurate prediction of protein secondary structural content. *Protein*  
491 *J.* **20**, 217–220 (2001).
- 492 57. Feng, Z. P. & Zhang, C. T. Prediction of membrane protein types based on the  
493 hydrophobic index of amino acids. *J. Protein Chem.* **19**, 269–275 (2000).
- 494 58. Sokal, R. R. & Thomson, B. A. Population structure inferred by local spatial  
495 autocorrelation: An example from an Amerindian tribal population. *Am. J. Phys.*  
496 *Anthropol.* **129**, 121–131 (2006).
- 497 59. Ong, S. A. K., Lin, H. H., Chen, Y. Z., Li, Z. R. & Cao, Z. Efficacy of different protein  
498 descriptors in predicting protein functional families. *BMC Bioinformatics* **8**, 1–14 (2007).
- 499 60. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.* **12**,  
500 2825–2830 (2012).
- 501 61. Chen, T. & Guestrin, C. XGBoost: A scalable tree boosting system. *Proc. ACM SIGKDD*  
502 *Int. Conf. Knowl. Discov. Data Min.* **13-17-Augu**, 785–794 (2016).
- 503 62. Edgar, R. C. MUSCLE: Multiple sequence alignment with high accuracy and high  
504 throughput. *Nucleic Acids Res.* (2004). doi:10.1093/nar/gkh340
- 505 63. Gouy, M., Guindon, S. & Gascuel, O. Sea view version 4: A multiplatform graphical user  
506 interface for sequence alignment and phylogenetic tree building. *Mol. Biol. Evol.* (2010).  
507 doi:10.1093/molbev/msp259
- 508 64. Lefort, V., Longueville, J. E. & Gascuel, O. SMS: Smart Model Selection in PhyML. *Mol.*  
509 *Biol. Evol.* (2017). doi:10.1093/molbev/msx149

- 510 65. Capra, J. A. & Singh, M. Predicting functionally important residues from sequence  
511 conservation. *Bioinformatics* (2007). doi:10.1093/bioinformatics/btm270
- 512 66. Greenbaum, J. *et al.* Functional classification of class II human leukocyte antigen (HLA)  
513 molecules reveals seven different supertypes and a surprising degree of repertoire sharing  
514 across supertypes. *Immunogenetics* **63**, 325–335 (2013).
- 515 67. Weiskopf, D. *et al.* Comprehensive analysis of dengue virus-specific responses supports  
516 an HLA-linked protective role for CD8+ T cells. *Proc. Natl. Acad. Sci. U. S. A.* **110**,  
517 E2046-53 (2013).
- 518 68. Jespersen, M. C., Peters, B., Nielsen, M. & Marcatili, P. BepiPred-2.0: Improving  
519 sequence-based B-cell epitope prediction using conformational epitopes. *Nucleic Acids*  
520 *Res.* **45**, W24–W29 (2017).
- 521 69. Schrödinger, L. The PyMol Molecular Graphics System, Versión 1.8. *Thomas Holder*  
522 (2015). doi:10.1007/s13398-014-0173-7.2
- 523 70. Zheng, W. *et al.* Deep-learning contact-map guided protein structure prediction in  
524 CASP13. *Proteins Struct. Funct. Bioinforma.* (2019). doi:10.1002/prot.25792
- 525
- 526

## 527 **Acknowledgments**

528 This work has been supported by the NIH-NIAID grant 1R01AI081062.

529

## 530 **Author contributions**

531 EO and YH contributed to the study design. EO, MW, AH collected the data. EO performed  
532 bioinformatics analysis. EO, MW, and YH wrote the manuscript. All authors performed result  
533 interpretation, and discussed and reviewed the manuscript.

534

535 **Competing financial interests:** The authors declare no competing financial interests.

## 536 **Figure Legends**

537

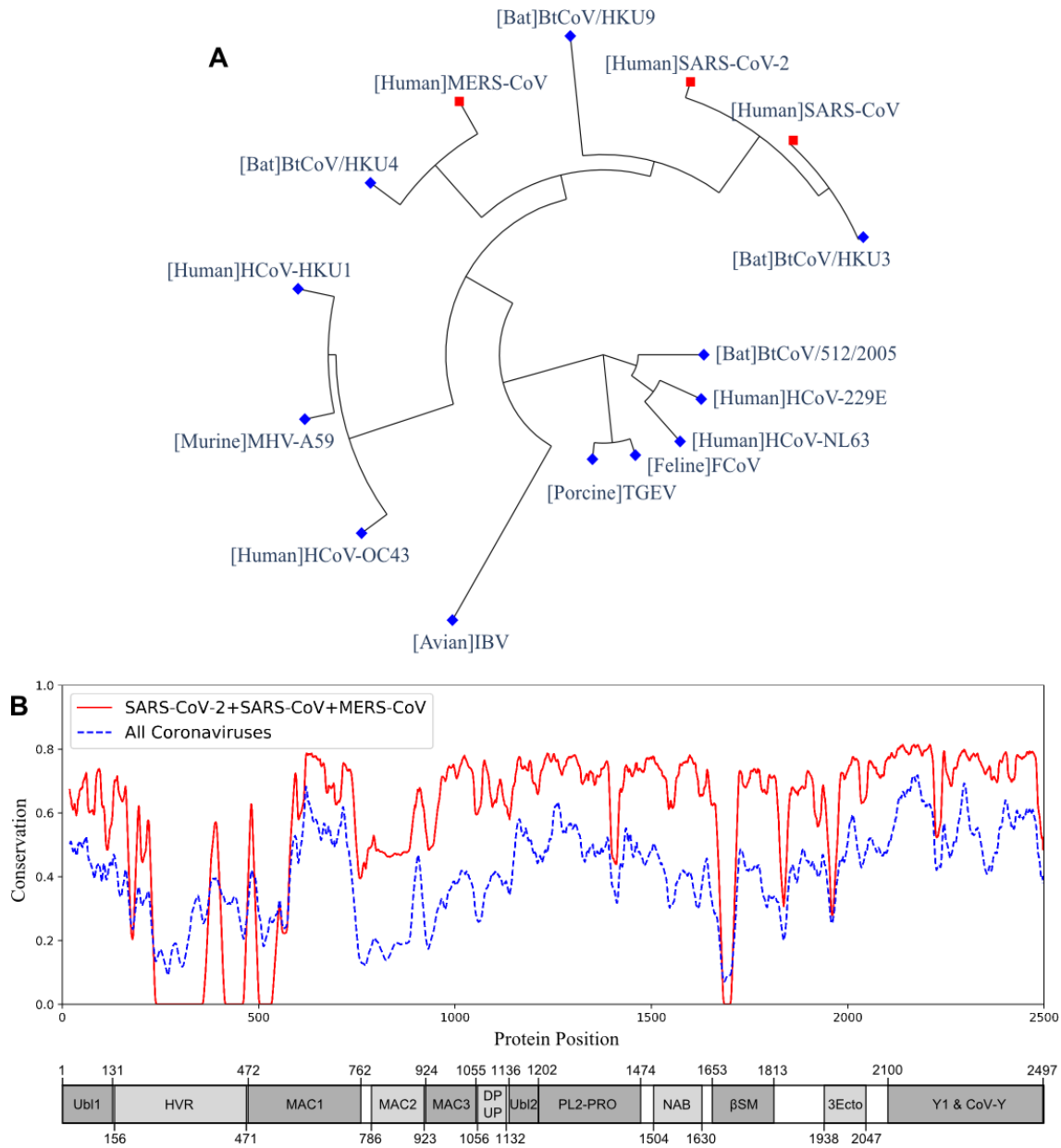
538 **Figure 1.** The phylogeny and sequence conservation of coronavirus nsp3. (A) Phylogeny of 15  
539 strains based on the nsp3 protein sequence alignment and phylogeny analysis. (B) The  
540 conservation of nsp3 among different coronavirus strains. The red line represents the  
541 conservation among the four strains (SARS-CoV, SARS-CoV-2, MERS, and BtCoV-HKU3).  
542 The blue line was generated using all the 15 strains. The bottom part represents the nsp3 peptides  
543 and their sizes. The phylogenetically close four strains have more conserved nsp3 sequences than  
544 all the strains being considered.

545

546 **Figure 2.** Predicted 3D structure of nsp3 protein highlighted with (A) MHC-I T cell epitopes  
547 (red), (B) MHC-II (blue) T cell epitopes, (C) linear B cell epitopes (green), and the merged  
548 epitopes. MHC-I epitopes are more internalized, MHC-II epitopes are more mixed, and B cells  
549 are more shown on the surface.

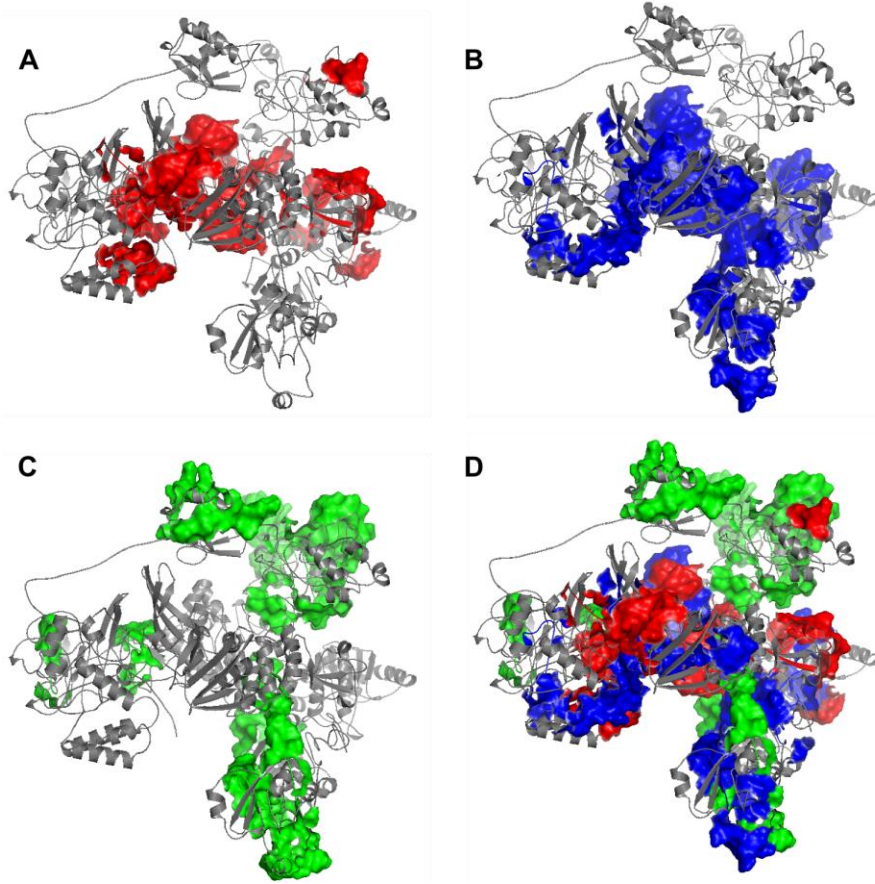
550

551 **Figure 3.** Immunogenic region of nsp3 between SARS-CoV-2 and the four conservation strains.  
552 (A) MHC-I (red) T cell epitope (B) MHC-II (blue) T cell epitope (C) linear B cell epitope  
553 (green).



554

555 **Figure 1.** The phylogeny and sequence conservation of coronavirus nsp3. (A) Phylogeny of 15  
 556 strains based on the nsp3 protein sequence alignment and phylogeny analysis. (B) The  
 557 conservation of nsp3 among different coronavirus strains. The red line represents the  
 558 conservation among the four strains (SARS-CoV, SARS-CoV-2, MERS, and BtCoV-HKU3).  
 559 The blue line was generated using all the 15 strains. The bottom part represents the nsp3 peptides  
 560 and their sizes. The phylogenetically close four strains have more conserved nsp3 sequences than  
 561 all the strains being considered.



562

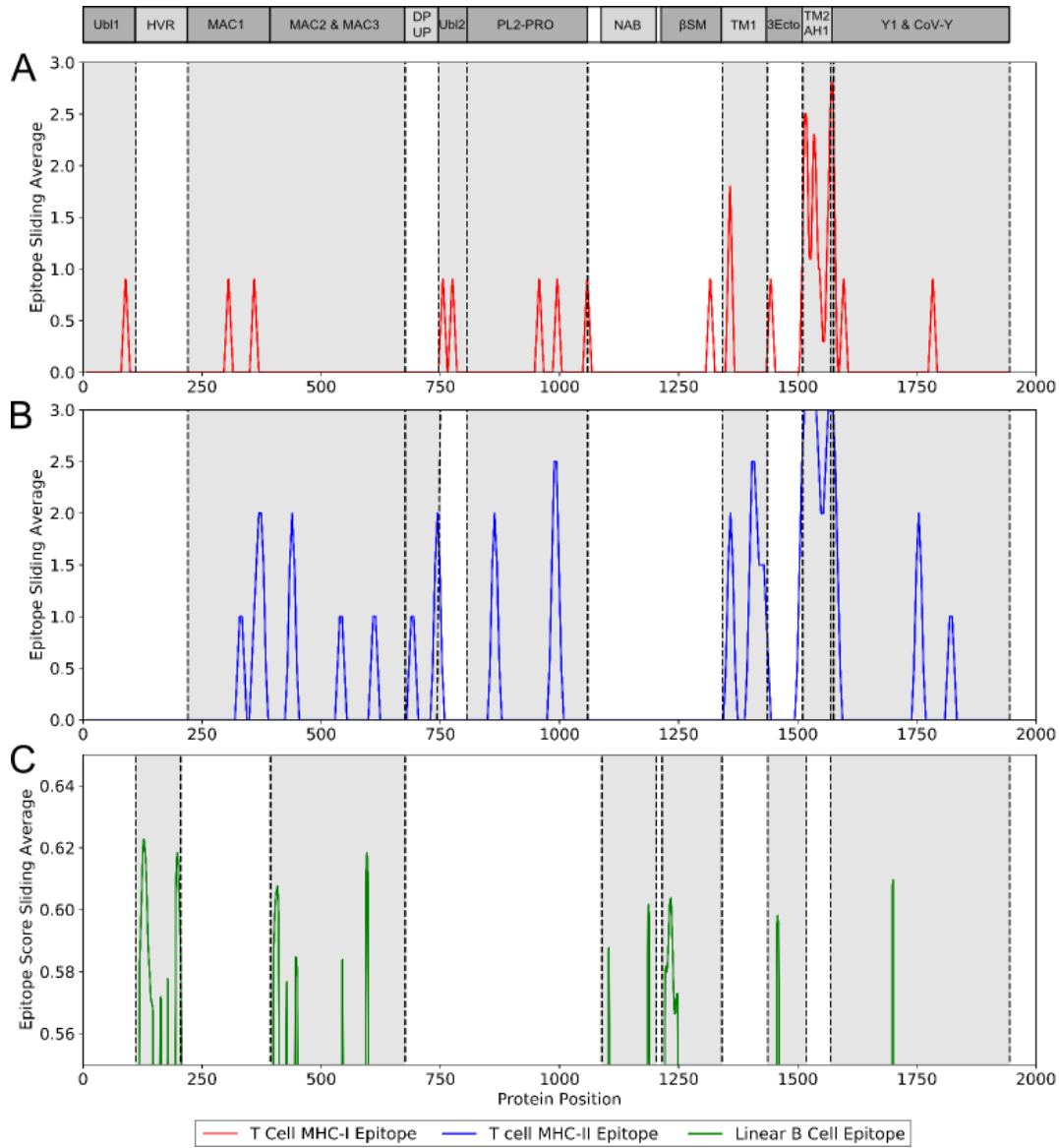
563

564 **Figure 2.** Predicted 3D structure of nsp3 protein highlighted with (A) MHC-I T cell epitopes

565 (red), (B) MHC-II (blue) T cell epitopes, (C) linear B cell epitopes (green), and the merged

566 epitopes. MHC-I epitopes are more internalized, MHC-II epitopes are more mixed, and B cells

567 are more shown on the surface.



568

569

570 **Figure 3.** Immunogenic region of nsp3 between SARS-CoV-2 and the four conservation strains.

571 (A) MHC-I (red) T cell epitope (B) MHC-II (blue) T cell epitope (C) linear B cell epitope

572 (green).

573 **Table 1.** Reported SARS-CoV, MERS-CoV, SARS-CoV-2 vaccine clinical trials.

<b>Virus</b>	<b>Location</b>	<b>Phase</b>	<b>Year</b>	<b>Identifier</b>	<b>Vaccine Type</b>
SARS-CoV	United States	I	2004	NCT00099463	recombinant DNA vaccine (S protein)
SARS-CoV	United States	I	2007	NCT00533741	whole virus vaccine
SARS-CoV	United States	I	2011	NCT01376765	recombinant protein vaccine (S protein)
MERS	United Kingdom	I	2018	NCT03399578	vector vaccine (S protein)
MERS	Germany	I	2018	NCT03615911	vector vaccine (S protein)
MERS	Saudi Arabia	I	2019	NCT04170829	vector vaccine (S protein)
MERS	Germany, Netherland	I	2019	NCT04119440	vector vaccine (S protein)
MERS	Russia	I,II	2019	NCT04128059	vector vaccine (protein not specified)
MERS	Russia	I,II	2019	NCT04130594	vector vaccine (protein not specified)
SARS-CoV2	United States	I	2020	NCT04283461	mRNA-based vaccine (S protein)
SARS-CoV2	China	I	2020	NCT04313127	vector vaccine (S protein)

574

575 **Table 2.** Vaccines tested for SARS-CoV and MERS-CoV.

Vaccine name	Vaccine type	Antigen	PMID
<b>SARS vaccines</b>			
CTLA4-S DNA vaccine	DNA	S	15993989
<i>Salmonella</i> -CTLA4-S DNA vaccine	DNA	S	15993989
<i>Salmonella</i> -tPA-S DNA vaccine	DNA	S	15993989
Recombinant spike polypeptide vaccine	Recombinant	S	15993989
N protein DNA vaccine	DNA	N	15582659
M protein DNA vaccine	DNA	M	16423399
N protein DNA vaccine	DNA	N	16423399
N+M protein DNA vaccine	DNA	N, M	16423399
tPA-S DNA vaccine	DNA	S	15993989
$\beta$ -propiolactone-inactivated SARS-CoV vaccine	Inactivated virus	whole virus	16476986
MA-ExoN vaccine	Live attenuated	MA-ExoN	23142821
rMA15- $\Delta$ E vaccine	Live attenuated	MA15	23576515
Ad S/N vaccine	Viral vector	S,N	16476986
ADS-MVA vaccine	Viral vector	S	15708987
MVA/S vaccine	Viral vector	S	15096611
<b>MERS vaccines</b>			
England1 S DNA Vaccine	DNA	S	26218507
MERS-CoV pcDNA3.1-S1 DNA vaccine	DNA	S	28314561
Inactivated whole MERS-CoV (IV) vaccine	Inactivated virus	whole virus	29618723
England1 S DNA +England1 S protein subunit Vaccine	Mixed	S1	26218507
England1 S1 protein subunit Vaccine	Subunit	S1	26218507
MERS-CoV S vaccine	Subunit	S	29618723
rNTD vaccine	Subunit	NTD of S	28536429
rRBD vaccine	Subunit	RBD of S	28536429
Ad5.MERS-S vaccine	Viral vector	S	25192975
Ad5.MERS-S1 vaccine	Viral vector	S1 subunit	25192975
VSV $\Delta$ G-MERS vaccine	Viral vector	S	29246504

576 Abbreviation: S, surface glycoprotein; N, nucleocapsid phosphoprotein; M, membrane glycoprotein; Exon,  
577 exoribonuclease; NTD, N-terminal domain; RBD, receptor binding domain.



578 **Table 3.** Vaxign-ML Prediction and adhesin probability of all SARS-CoV-2 proteins.

	<b>Protein</b>	<b>Vaxign-ML Score</b>	<b>Adhesin Probability</b>	
orf1ab	nsp1	Host translation inhibitor	79.312	
	nsp2	Non-structural protein 2	89.647	
	nsp3	Non-structural protein 3	<b>95.283*</b>	<b>0.524<sup>#</sup></b>
	nsp4	Non-structural protein 4	89.647	0.289
	3CL-PRO	Proteinase 3CL-PRO	89.647	<b>0.653<sup>#</sup></b>
	nsp6	Non-structural protein 6	89.017	0.320
	nsp7	Non-structural protein 7	89.647	0.269
	nsp8	Non-structural protein 8	<b>90.349*</b>	<b>0.764<sup>#</sup></b>
	nsp9	Non-structural protein 9	89.647	<b>0.796<sup>#</sup></b>
	nsp10	Non-structural protein 10	89.647	<b>0.769<sup>#</sup></b>
	RdRp	RNA-directed RNA polymerase	89.647	0.229
	Hel	Helicase	89.647	0.398
	ExoN	Guanine-N7 methyltransferase	89.629	0.183
	NendoU	Uridylate-specific endoribonuclease	89.647	0.254
	2'-O-MT	2'-O-methyltransferase	89.647	0.421
	S	Surface glycoprotein	<b>97.623*</b>	<b>0.635<sup>#</sup></b>
	ORF3a	ORF3a	66.925	0.383
	E	envelope protein	23.839	0.234
	M	membrane glycoprotein	84.102	0.282
	ORF6	ORF6	33.165	0.095
ORF7	ORF7a	11.199	0.451	
ORF8	ORF8	31.023	0.311	
N	nucleocapsid phosphoprotein	89.647	0.373	
ORF10	ORF10	6.266	0.0	

579 \* denotes Vaxign-ML predicted vaccine candidate.

580 # denotes predicted adhesin.

581

582