1

2

# Novel Immunoglobulin Domain Proteins Provide Insights into Evolution and Pathogenesis Mechanisms of SARS-Related Coronaviruses

5

6

Yongjun Tan[a], Theresa Schneider[a], Matthew Leong[b], L Aravind[c], Dapeng Zhang[a,d,*]

8

9

10

[a] Department of Biology, College of Arts and Sciences, Saint Louis University, MO 63110

[b] School of Medicine, Saint Louis University, MO 63110

[c] National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD 20894

[d] Program of Bioinformatics and Computational Biology, College of Arts and Sciences, Saint Louis University, MO 63110

*Corresponding Author: Dapeng Zhang, Ph.D. (dapeng.zhang@slu.edu)

18

19 **ABSTRACT**

20 A novel coronavirus (SARS-CoV-2) is the causative agent of an emergent severe respiratory disease

21 (COVID-19) in humans that is threatening to result in a global health crisis. By using genomic, sequence,

22 structural and evolutionary analysis, we show that Alpha- and Beta-CoVs possess several novel families

23 of immunoglobulin (Ig) domain proteins, including ORF8 and ORF7a from SARS-related coronaviruses

24 and two protein groups from certain Alpha-CoVs. Among them, ORF8 is distinguished in being rapidly

25 evolving, possessing a unique insert and a hypervariable position among SARS-CoV-2 genomes in its

26 predicted ligand-binding groove. We also uncover many Ig proteins from several metazoan viruses

27 which are distinct in sequence and structure but share an architecture comparable to that of CoV Ig

28 domain proteins. Hence, we propose that deployment of Ig domain proteins is a widely-used strategy by

29 viruses, and SARS-CoV-2 ORF8 is a potential pathogenicity factor which evolves rapidly to counter the

30 immune response and facilitate the transmission between hosts.

31

32 **KEYWORDS** Coronavirus, COVID-19, SARS, ORF8, Immunoglobin, Evolution, Pathogenicity, Immune

33 evasion

**Introduction**

34

35    Nidoviruses are an ancient group of lipid-enveloped viruses with non-segmented RNA genomes, which

36    are known to infect oomycetes and animals, including molluscs, arthropods and vertebrates (*1*). Among

37    them are the coronaviruses (CoVs) which possess the largest known monopartite RNA genomes and are

38    classified into four genera—Alphacoronavirus, Betacoronavirus, Gammacoronavirus, and

39    Deltacoronavirus (*2*).  Over the past two decades, Beta-CoVs, including the viruses responsible for

40    Severe Acute Respiratory Syndrome (SARS) in 2003 and Middle Eastern Respiratory Syndrome (MERS) in

41    2012, have emerged as significant local land global health concerns with economic consequences (*3, 4*).

42    Recently, a novel severe respiratory disease has emerged in humans (abbreviated COVID-19) (*5, 6*).

43    COVID-19 presents with a relatively long incubation period of 1-2 weeks followed by development of

44    fever, dry cough, dyspnea and bilateral ground-glass opacities in the lungs (*7*). In some patients, this can

45    proceed to fatal respiratory failure, characterized by acute lung injury (*8*) and acute respiratory distress

46    syndrome (*9*). Within several months of the first outbreak, there have been over 75,000 confirmed cases

47    with over 2,000 deaths of COVID-19 globally (https://www.who.int/docs/default-

48    source/coronavirus/situation-reports/20200220-sitrep-31-covid-19.pdf). Due to the rapid spread and

49    potential severity of the disease, COVID-19 poses as a potential major threat to human health. A novel

50    coronavirus (SARS-CoV-2) was identified as the causative agent of COVID-19, and phylogenomic analysis

51    has shown that it belongs to the same larger clade of Beta-CoVs as the SARS-CoV with a likely origin in

52    bats (*5, 10, 11*). Despite intense scrutiny, multiple proteins encoded by the SARS-related CoV (including

53    SARS-CoV-2) genome remain enigmatic. Here, we present a computational and evolutionary analysis to

54    show that one such mysterious protein, ORF8, and several others from Alpha- and Beta-coronavirus,

55    comprise novel families of immunoglobulin domain proteins, which might function as potential

56    modulators of host immunity to delay or attenuate the immune response against the viruses.

57

58    **Materials and methods**

59    ***Genome comparison analysis***

60    We retrieved the SARS-related CoV genomes by searching against the non-redundant (nr) nucleotide

61    database of the National Center for Biotechnology Information (NCBI) with the SARS-CoV-2 genome

62    sequence (NC_045512.2) as a query (*12*). The program CD-HIT was used for similarity-based clustering

63    (*13*). A multiple sequence alignment of whole virus genomes was performed by KALIGN (*14*). Based on

64    the MSA, a similarity plot was constructed by a custom Python script, which calculated the identity

65    between each subject sequence and the SARS-CoV-2 genome sequence based on a custom sliding

66  window size and step size. Open reading frames of virus genomes used in this study were extracted

67  from an NCBI Genbank file.

68

69  *Protein sequence analysis*

70  To collect protein homologs, iterative sequence profile searches were conducted by the programs PSI-

71  BLAST (Position-Specific Iterated BLAST)(*12*) and JACKHMMER (*15*), which searched against the non-

72  redundant (nr) protein database of NCBI with a cut-off e-value of 0.005 serving as the significance

73  threshold. Similarity-based clustering was conducted by BLASTCLUST, a BLAST score-based single-linkage

74  clustering method (ftp://ftp.ncbi.nih.gov/blast/documents/blastclust.html). Multiple sequence

75  alignments were built by the KALIGN (*14*), MUSCLE(*16*) and PROMALS3D(*17*) programs, followed by

76  careful manual adjustments based on the profile–profile alignment, the secondary structure information

77  and the structural alignment. Profile-profile comparison was conducted using the HHpred program (*18*).

78  The consensus of the alignment was calculated using a custom Perl script. The alignments were colored

79  using an in-house alignment visualization program written in perl and further modified using adobe

80  illustrator. Signal peptides were predicted by the SignalP-5.0 Server (*19*). The transmembrane regions

81  were predicted by the TMHMM Server v. 2.0 (*20*).

82

83  *Identification of distinct viral Ig domain proteins*

84  By using the protein remote relationship detection methods, we generated a collection of distinct Ig

85  domains from the Pfam database (*21*) and also from our local domain database. Then, we utilized the

86  hmmscan program of the HMMER package (*22*) and RPS-BLAST (*12, 23*) to retrieve the homologs from

87  viral genomes.

88

89  *Molecular Phylogenetic analysis*

90  The evolutionary history was inferred by using the Maximum Likelihood method based on the JTT

91  w/freq. model (*24*). The tree with the highest log likelihood is shown. Support values out of 100

92  bootstraps are shown next to the branches (*25*). Initial tree(s) for the heuristic search were obtained

93  automatically by applying Neighbor-Join and BioNJ algorithms to a matrix of pairwise distances

94  estimated using a JTT model, and then selecting the topology with the superior log likelihood value. A

95  discrete Gamma distribution was used to model evolutionary rate differences among sites (4

96  categories). The rate variation model allowed for some sites to be evolutionarily invariable. The tree is

97    drawn to scale, with branch lengths measured in the number of substitutions per site. The tree diagram

98    was generated using MEGA Tree Explorer (*26*)

99

100   ***Entropy analysis***

101   Position-wise Shannon entropy (H) for a given multiple sequence alignment was calculated using the

102   equation:

$$H = -\sum_{i=1}^{M} P_i \log_2 P_i$$

103

104   *P* is the fraction of residues of amino acid type *i*, and *M* is the number of amino acid types. The Shannon

105   entropy for the *i*th position in the alignment ranges from 0 (only one residue at that position) to 4.32 (all

106   20 residues equally represented at that position). Analysis of the entropy values which were thus

107   derived was performed using the R language.

108

109   ***Protein Structure Prediction and Analysis***

110   The secondary structural prediction was conducted using the Jnet (Joint Network) program (*27*). Jnet is a

111   neural network-based predictor which trains neural networks from three different types of profiles:

112   profile PSSM, profile HMM, and residue frequency profile. It generates a consensus secondary structure

113   with an average accuracy of 72% or greater.

114

115   The Modeller9v1 program (Sali and Blundell, 1993) was utilized for homology modeling of structure of

116   SARS-CoV-2 ORF8 using the SARS-CoV ORF7 (1xak_A) (*28*) as a template. Since in these low sequence-

117   identity cases, sequence alignment is the most important factor affecting the quality of the model

118   (Cozzetto and Tramontano, 2005), alignments used in this study have been carefully built and cross-

119   validated based on the information from HHpred and edited manually using the secondary structure

120   information. We generated five models and selected the one that had better model accuracy p-value

121   and global model quality score as assessed by ModFOLD6 online server (*29*). Structural analysis and

122   comparison were conducted using the molecular visualization program PyMOL (*30*). The structural

123   similarity search was performed using the DALI server (*31*).

124

125   **Results and discussion**

126   ***Comparative genomics unveils fast-evolving regions of SARS-related CoV genomes***

127    The host-pathogen arms-race has selected for a disparate complement of viral genes involved in

128    pathogenesis. These genes often rapidly diversify through recombination and mutations to keep up with

129    the evolution of host resistance. To identify proteins with potential pathogenic roles in COVID-19, we

130    conducted a comparative genomic analysis of the coronaviruses. Similarity plots show that the bat CoV

131    RaTG13 is the closest relative of SARS-CoV-2 with no evidence for recombination between them (Figure

132    1 and Figure S1). SARS-CoV-2 also shows good similarity to two other bat viruses, CoVZXC21 and

133    CoVZXC45, first in the 5' half of ORF1 and again after nucleotide number 20,000 of the genome.

134    However, the remaining part of ORF1 of SARS-CoV-2 and RatG13 show no specific relationship to these

135    viruses. This suggests a recombination event between the common ancestor of SARS-CoV-2 and RaTG13,

136    and probably another member of the SARS-related clade close to CoVZXC21 and CoVZXC45. In addition

137    to this major recombination event, we identified several smaller regions which might have undergone

138    recombinational diversification during the emergence of the SARS-CoV-2 genome (Figure S1). Notably,

139    many of them are clustered in three regions displaying extensive diversification, corresponding to the N-

140    terminal region of the ORF1a polyprotein, the Spike protein, and the uncharacterized protein encoded

141    by ORF8 (Figure 1 and Figure S1).

142

143    ***Identification of novel immunoglobin protein families in CoVs***

144    Among these three fast-evolving proteins, we focused on the ORF8 protein as it is one of the so-called

145    accessory proteins, which does not participate in viral replication (*32, 33*), raising the possibility that it

146    might have a role in viral pathogenesis. It is a predicted secreted protein present only in some Beta-

147    CoVs, including SARS-CoV-2 but not the MERS-like clade. Profile-profile comparisons using a sequence-

148    profile built from the multiple sequence alignment of all available ORF8 proteins showed it to be

149    unexpectedly homologous to the membrane-anchored ORF7a protein from the same subset of Beta-

150    CoVs, and several proteins (variously annotated as ORF9 or ORF10) from a subset of bat Alpha-CoVs

151    (Figure 2A) (probability=94% of profile-profile match) (*34*). ORF7a is a known member of the

152    immunoglobulin (Ig) domain superfamily and is specifically related to extracellular metazoan Ig domains

153    that are involved in adhesion, such as ICAM (*35, 36*). The Beta-CoV ORF8, ORF7a and the Alpha-CoV Ig

154    domains display a classic β-sandwich fold with seven β-stands and share the characteristic pattern of

155    two cysteines which form stabilizing disulfide bonds with metazoan Ig domains (Figure 2A and Figure S2)

156    (*37*). However, they are unified as a clade by the presence of an additional pair of conserved disulfide-

157    bonding cysteines (Figure 2A). Nonetheless, there are notable structural differences between the three

158    groups of proteins. ORF8 is distinguished from ORF7a and Alpha-CoV Ig proteins by the loss of the C-

159    terminal transmembrane (TM) helix and the acquisition of a long insert between stands 3 and 4 with a

160    conserved cysteine which might facilitate dimerization through disulfide-bond formation (Figure 2A).

161    The homology model of ORF8, based on the structure of SARS-CoV ORF7a (pdb id:1XAK), suggests that

162    this insert augments a potential peptide-ligand binding groove that has been proposed for ORF7a

163    (Figure 2B and Figure S3). Hence, the emergence of the insert has gone hand-in-hand with the

164    acquisition of a modified interaction interface.

165

166    Besides these families, we identified a fourth family of Ig domains from the same Alpha-CoVs which

167    contain the above-discussed ORF7a/8-like Ig family (Figure 3A). These Alpha-CoVs typically possess one

168    or two paralogous copies annotated as either ORF4a/b or NS5a/b. From their sequences, these Ig

169    domains are not closely related to the ORF7a and ORF8 Ig domains (Figure S4). However, profile-profile

170    searches have shown that they are related to Ig domains found in the adenoviral E3-CR1 proteins

171    (probability: 90% of matching the Pfam CR1 Ig domain profile) (Figure S4). In these searches, they also

172    yield weaker hits to two other Ig domains, namely the poxviral decoy interferon receptors and human T-

173    cell surface CD3 zeta (Figure S4).

174

175    ***ORF8 is a fast-evolving protein in SARS-related CoVs***

176    Phylogenetic analysis of ORF7a, ORF8 and Alpha-CoV Ig domains shows that each group represents a

177    distinct clade (Figure 2C). The tree topology of ORF7a mirrors that of the polymerase tree (Figure 3A);

178    however, the topology of the ORF8-Ig clade is not consistent with it. This might be due to a

179    recombination event between the SARS-related CoVs (as suggested by the similarity plot analysis)

180    and/or unusual divergence under selection. To better understand the functional difference between

181    ORF8 and ORF7a, we examined the column-wise Shannon entropy in the 20 amino acid alphabet and

182    found that the ORF8 has significantly higher mean entropy than ORF7a (ORF8: 1.09 vs ORF7a: 0.22, p<

183    $10^{-16}$ for the $H_0$ of congruent means by t-test) (Figure 2D). By comparing column-wise entropies in both

184    the 20 amino acid and a reduced 8-letter alphabet (where amino acids are grouped based on similar

185    side-chain chemistry), we found at least 14 positions in ORF8 which show high entropy in both alphabets

186    as compared to a single position in ORF7a (Figure 2D). This indicates that ORF8 is a fast-evolving protein

187    under selection for diversity as contrast to ORF7a. Strikingly, one of these highly variable positions,

188    which features residues with very different side characters (hydrophobic, acidic, alcoholic and proline),

189    corresponding to Leu84 was also identified as the most variable position across 54 closely related

190    human SARS-CoV-2 genome sequences (*38*). In our structural model, this residue is positioned at the

191    predicted peptide-ligand binding groove of the ORF8-Ig domain (Figure 2B). Therefore, our entropy and

192    structural analysis of the ORF8-Ig domain, in conjunction with its hypervariable position found in human

193    SARS-CoV-2 genomes, points to a role for ORF8 at the interface of the host-virus interaction possibly in a

194    pathogenic context.

195

196    ***Ig domain proteins are newly acquired in subsets of Alpha- and Beta-CoVs***

197    We examined the distribution of CoV Ig proteins in the context of a phylogenetic tree of both Beta-and

198    Alpha-CoVs based on their polymerase proteins (Figure 3A). Other than the two subsets of Beta-CoVs

199    and Alpha-CoVs that contain the above-described Ig domain proteins, no other CoVs contain any Ig

200    domain proteins (Figure 3A). The immediate sister-groups of the Ig-containing CoVs typically have Spike,

201    E, M and N, and one or two other uncharacterized accessory proteins which are not Ig domains to our

202    best knowledge. The Alpha-CoV ORF9/10 share a C-terminal TM helix and, along with ORF7a of the Beta-

203    CoVs, lack the insert in the Ig domain (Figure 2A). Hence, it is possible that this architecture represents

204    the ancestral state which was present in the common ancestor of both Alpha-CoVs and Beta-CoVs.

205    Under this scenario, the protein was displaced/lost both in certain Alpha-CoVs and Beta-CoVs.

206    Alternatively, ORF7a could have been exchanged between Alpha- and Beta-CoVs by a recombination

207    event. In both scenarios, ORF8 arose likely via a duplication of ORF7a in specifically the Beta-CoVs.

208    Although we couldn't identify the ultimate precursors of the CoV Ig domains, they are likely to have

209    been acquired on at least two independent occasions from different sources. The CoV ORF7a-ORF8

210    families might have derived from the metazoan adhesion Ig families, and the ORF4a/b-like Ig domains of

211    Alpha-CoVs were likely acquired from adenoviral CR1 Ig domains with which they share some specific

212    sequence features.

213

214    ***Divergent Ig proteins with a comparable architecture are deployed by distinct viruses***

215    The presence of multiple Ig domains with different affinities in CoVs prompted us to more generally

216    survey animal viruses for Ig domains. By using the Pfam models (*39*) and our own PSSMs created from

217    PSI-BLAST runs with Ig domains (*12*), we were able to identify about 17 distinct viral Ig domain families

218    in a wider diversity of animal viruses (Figure 3B and Table S2). In addition to CoVs, such Ig domain

219    proteins can be found in adenoviruses, NCLDVs, Herpesviruses and Phenuiviruses. These viral Ig domains

220    are highly divergent; many of them are only found in certain viral groups. However, the majority have an

221    architecture comparable to the CoV-Ig domains, with an N-terminal signal peptide, one or multiple Ig

222    domains and a C-terminal TM region often followed by a stretch of basic residues. Thus, although the Ig

223   domains are not the universally preserved component of viruses, they have been acquired and selected

224   independently by a wide range of viruses. The presence of a proofreading 3'-5' exonuclease has been

225   proposed to favor the emergence of larger RNA genomes in CoVs (*40*). Indeed, this might have also

226   contributed to the acquisition of potential pathogenesis factors such as the Ig domains described herein

227   which are comparable to those seen in DNA viruses.

228

229   ***Novel CoV Ig domain proteins are potential immune modulators***

230   Why have diverse viruses independently acquired the Ig domain during their evolution? First, the Ig

231   domains are major mediators of adhesive interactions in both eukaryotes and prokaryotes (*37, 41, 42*).

232   Thus, this domain can be used for adherence for cell to cell spread (e.g. herpesviral Ig domain proteins)

233   (*43*). Further, Ig domains are major building blocks of metazoan immune systems. Thus, viruses often

234   utlize this domain to disrupt immune signaling of the host. For example, in adenoviruses, the CR1 Ig

235   domain proteins have been shown to block the surface expression on infected cells of class I major

236   histocompatibility complex molecules by blocking their trafficking from the endoplasmic reticulum (ER)

237   to Golgi (*44*). This has been shown to affect the host inflammatory response and modulates the

238   presentation of viral antigens to T-cells (*45*). In poxvirus, the secreted Ig domain proteins function as

239   interferon receptors or decoys that bind the interferon-α/β and disrupt signaling via the endogenous

240   host receptors (*46*). Further, SARS-ORF7a has been implicated in the interaction with bone marrow

241   stromal antigen 2 (BST-2), which tethers budding virions to the host cell in a broad-spectrum antiviral

242   response, to prevent the N-linked glycosylation of BST-2 thereby crippling the host response against the

243   virus (*47*).  Given their shared evolutionary history and similar sequence and structural features, we

244   predict that the newly identified CoV Ig domain proteins, such as ORF8 of SARS-CoV-2, might similarly

245   function as immune modulators.

246

247   While ORF8 is a paralog of ORF7a, its lack of the TM segment, unique insert and significantly more rapid

248   evolution suggest that it has acquired a distinct function and has been under strong positive selection

249   One possible mechanism is that, like the adenoviral CR1 proteins, it interferes with MHC molecules to

250   attenuate antigen presentation, resulting in ineffective detection of the virus by the host immune

251   system. Consistent with this, the SARS-CoV ortholog translocates to the ER (*48*) and its higher variability

252   indicates probable selection due to its interaction with a rapidly evolving host molecule. Notably, while

253   the SARS-CoV ORF8 isolates from civets and early stages of the human epidemic are intact, it split up

254   into two ORFs (ORF8a and ORF8b) during the subsequent human epidemic (*49*). ORF8a and ORF8b

255     retain the conserved Cys residues of the Ig domain and have been observed to form a complex in a

256     yeast-two hybrid interaction study (*50*). This suggests it might still fold into an intact structure held by

257     the two disulfide bonds formed by four conserved Cys residues.

258

259     In conclusion, the presence of fast-evolving ORF8 Ig domain proteins in the SARS-related viruses,

260     including the emergent 2019 SARS-CoV-2, suggests that they might be potential pathogenicity factors

261     which counter or attenuate the host immune response and might have facilitated the transmission

262     between hosts. We hope that the discovery and analyses of the novel Ig domain proteins reported here

263     will help the community better understand the evolution and pathogenesis mechanisms of these

264     coronaviruses.

265

266

**Acknowledgments**

270

**Disclosure statement**

The authors declare no competing interests.

273

274

275    **Figures**

276    **Figure 1. Genome comparison analysis of SARS-related CoVs.**

277    Similarity plot of SARS-related CoVs against human SARS-CoV-2 Wuhan-Hu-1 genome (NC_045512.2),

278    based on a multiple sequence alignment of the whole genomes. Each point represents the percent

279    identity of a 200 bp window of the alignment with a 50 bp step size between each point. The open

280    reading frames of the SARS-CoV-2 genome (NC_045512.2) are shown above the plot. Each colored line

281    corresponds to the nucleotide similarity between the human SARS-CoV-2 genome (NC_045512.2) and

282    the respective SARS-related CoV genome. The red arrows and dashed line surround a region displaying

283    major divergence due to possible recombination within SARS-related CoV genomes. The regions marked

284    by a solid red line highlight fast-evolving regions among the SARS-related CoV genomes. For detailed

285    information about the genomes that were used in this study, refer to Table S1.



286

287  **Figure 2. Sequence, structure and evolutionary analysis of novel Ig domain proteins in SARS-related**

288  **CoVs.**

289  (A) Multiple sequence alignment (MSA) and representative domain architectures of ORF7a-Ig, ORF8-Ig,

290  and ORF7a/8-like Ig domain families. Each sequence in the MSA was labelled by its species abbreviation

291  followed by its source. The predicted secondary structure is shown above each alignment and the

292  consensus is shown below the super-alignment, where h stands for hydrophobic residues, s for small

293  residues, and p for polar residues. Two pairs of conserved cysteines that form disulfide bonds are

294  highlighted in red. (B) Homology model of SARS-CoV-2 ORF8-Ig domain (YP_009724396.1) and the

295  location of the hypervariable position corresponding to Leu84 in the predicted ligand-binding groove.

296  The β-sheets of the common core of the Ig fold are colored in blue, the insert in ORF8-Ig in orange and

297  the loops in grey. The characteristic disulfide bonds are highlighted in yellow. (C) Maximum likelihood

298  phylogenetic analysis of CoV Ig domain families. Support values out of 100 bootstraps are shown for the

299  major branches only. (D) Entropy plot for the ORF7a and ORF8 proteins in betacoronavirus. Left:

300  Shannon entropy computed for each column for a character space of 20 amino acids and presented as

301  mean entropy in a sliding window of 30 residues. The mean entropy across the entire length of the

302  protein is indicated as a green horizontal line. Right: Shannon entropy in regular amino acid alphabet (20

303  amino acids) are shown above the zero line in shades of orange. Shannon entropy in a reduced alphabet

304  of 8 residues are shown below the zero line in shades of blue. If a position shows high entropy in both

305  alphabets it is a sign of potential positive selection at those positions for amino acids of different

306  chemical character.

307

308

309    **Figure 3.** (A) Genomic structure analysis of SARS-related CoVs. The tree of coronavirus was built based

310    on an MSA of a coronavirus RNA-directed, RNA polymerase domain using a maximum likelihood model.

311    Supporting values from 100 bootstraps are shown for the major branches only. The genome structure of

312    major representative CoVs are shown right to the terminal clade of the phylogenetic tree. (B)

313    Representative domain architectures of the Ig components in different animal viruses. Proteins were

314    grouped based on their families, except for proteins of coronavirus, which were grouped based on their

315    genus. For the information of the NCBI accession numbers, refer to the supplementary data.

**References**

1.  S. Perlman, T. Gallagher, E. J. Snijder, *Nidoviruses*. (ASM Press, Washington, DC, 2008), pp. xvi, 433 p.

2.  J. Cui, F. Li, Z. L. Shi, Origin and evolution of pathogenic coronaviruses. *Nat Rev Microbiol* **17**, 181-192 (2019).

3.  M. A. Marra *et al.*, The Genome sequence of the SARS-associated coronavirus. *Science* **300**, 1399-1404 (2003).

4.  A. M. Zaki, S. van Boheemen, T. M. Bestebroer, A. D. Osterhaus, R. A. Fouchier, Isolation of a novel coronavirus from a man with pneumonia in Saudi Arabia. *N Engl J Med* **367**, 1814-1820 (2012).

5.  P. Zhou *et al.*, A pneumonia outbreak associated with a new coronavirus of probable bat origin. *Nature*, (2020).

6.  N. Zhu *et al.*, A Novel Coronavirus from Patients with Pneumonia in China, 2019. *N Engl J Med*, (2020).

7.  N. Chen *et al.*, Epidemiological and clinical characteristics of 99 cases of 2019 novel coronavirus pneumonia in Wuhan, China: a descriptive study. *Lancet* **395**, 507-513 (2020).

8.  J. P. Kanne, Chest CT Findings in 2019 Novel Coronavirus (2019-nCoV) Infections from Wuhan, China: Key Points for the Radiologist. *Radiology*, 200241 (2020).

9.  C. Huang *et al.*, Clinical features of patients infected with 2019 novel coronavirus in Wuhan, China. *Lancet* **395**, 497-506 (2020).

10. D. Paraskevis *et al.*, Full-genome evolutionary analysis of the novel corona virus (2019-nCoV) rejects the hypothesis of emergence as a result of a recent recombination event. *Infect Genet Evol* **79**, 104212 (2020).

11. L. Zhang, F. M. Shen, F. Chen, Z. Lin, Origin and evolution of the 2019 novel coronavirus. *Clin Infect Dis*, (2020).

12. S. F. Altschul *et al.*, Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic acids research* **25**, 3389-3402 (1997).

13. W. Li, A. Godzik, Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics* **22**, 1658-1659 (2006).

14. T. Lassmann, E. L. Sonnhammer, Kalign–an accurate and fast multiple sequence alignment algorithm. *BMC bioinformatics* **6**, 298 (2005).

15. S. R. Eddy, in *Genome Informatics 2009: Genome Informatics Series Vol. 23*. (World Scientific, 2009), pp. 205-211.

16. R. C. Edgar, MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic acids research* **32**, 1792-1797 (2004).

17. J. Pei, B.-H. Kim, N. V. Grishin, PROMALS3D: a tool for multiple protein sequence and structure alignments. *Nucleic acids research* **36**, 2295-2300 (2008).

18. J. Söding, A. Biegert, A. N. Lupas, The HHpred interactive server for protein homology detection and structure prediction. *Nucleic acids research* **33**, W244-W248 (2005).

19. L. Holm, S. Kääriäinen, P. Rosenström, A. Schenkel, Searching protein structure databases with DaliLite v. 3. *Bioinformatics* **24**, 2780-2781 (2008).

20. J. J. A. Armenteros *et al.*, SignalP 5.0 improves signal peptide predictions using deep neural networks. *Nature biotechnology* **37**, 420-423 (2019).

21. A. Krogh, B. Larsson, G. Von Heijne, E. L. Sonnhammer, Predicting transmembrane protein topology with a hidden Markov model: application to complete genomes. *Journal of molecular biology* **305**, 567-580 (2001).

363    22.    S. El-Gebali *et al.*, The Pfam protein families database in 2019. *Nucleic acids research* **47**, D427-
364           D432 (2019).
365    23.    A. Krogh, M. Brown, I. S. Mian, K. Sjolander, D. Haussler, Hidden Markov models in
366           computational biology. Applications to protein modeling. *Journal of molecular biology* **235**,
367           1501-1531 (1994).
368    24.    A. Marchler-Bauer *et al.*, CDD: a database of conserved domain alignments with links to domain
369           three-dimensional structure. *Nucleic acids research* **30**, 281-283 (2002).
370    25.    D. T. Jones, W. R. Taylor, J. M. Thornton, The rapid generation of mutation data matrices from
371           protein sequences. *Bioinformatics* **8**, 275-282 (1992).
372    26.    J. Felsenstein, Confidence limits on phylogenies: an approach using the bootstrap. *Evolution* **39**,
373           783-791 (1985).
374    27.    S. Kumar, G. Stecher, K. Tamura, MEGA7: molecular evolutionary genetics analysis version 7.0
375           for bigger datasets. *Molecular biology and evolution* **33**, 1870-1874 (2016).
376    28.    A. Drozdetskiy, C. Cole, J. Procter, G. J. Barton, JPred4: a protein secondary structure prediction
377           server. *Nucleic acids research* **43**, W389-W394 (2015).
378    29.    N. Eswar *et al.*, Comparative protein structure modeling using Modeller. *Current protocols in
379           bioinformatics* **15**, 5.6. 1-5.6. 30 (2006).
380    30.    A. H. Maghrabi, L. J. McGuffin, ModFOLD6: an accurate web server for the global and local
381           quality estimation of 3D protein models. *Nucleic acids research* **45**, W416-W421 (2017).
382    31.    W. L. DeLano, Pymol: An open-source molecular graphics tool. *CCP4 Newsletter On Protein
383           Crystallography* **40**, 82-92 (2002).
384    32.    M. L. Dediego *et al.*, Pathogenicity of severe acute respiratory coronavirus deletion mutants in
385           hACE-2 transgenic mice. *Virology* **376**, 379-389 (2008).
386    33.    B. Yount *et al.*, Severe acute respiratory syndrome coronavirus group-specific open reading
387           frames encode nonessential functions for replication in cell cultures and mice. *J Virol* **79**, 14909-
388           14922 (2005).
389    34.    J. Soding, A. Biegert, A. N. Lupas, The HHpred interactive server for protein homology detection
390           and structure prediction. *Nucleic Acids Res* **33**, W244-248 (2005).
391    35.    C. A. Nelson, A. Pekosz, C. A. Lee, M. S. Diamond, D. H. Fremont, Structure and intracellular
392           targeting of the SARS-coronavirus Orf7a accessory protein. *Structure* **13**, 75-85 (2005).
393    36.    K. Hanel, T. Stangler, M. Stoldt, D. Willbold, Solution structure of the X4 protein coded by the
394           SARS related coronavirus reveals an immunoglobulin like fold and suggests a binding activity to
395           integrin I domains. *J Biomed Sci* **13**, 281-293 (2006).
396    37.    J. M. Berg, J. L. Tymoczko, L. Stryer, L. Stryer, *Biochemistry*.  (W.H. Freeman, New York, ed. 5th,
397           2002).
398    38.    C. Ceraolo, F. M. Giorgi, Genomic variance of the 2019-nCoV coronavirus. *J Med Virol*,  (2020).
399    39.    S. El-Gebali *et al.*, The Pfam protein families database in 2019. *Nucleic Acids Res* **47**, D427-D432
400           (2019).
401    40.    F. Ferron *et al.*, Structural and molecular basis of mismatch correction and ribavirin excision
402           from coronavirus RNA. *Proc Natl Acad Sci U S A* **115**, E162-E171 (2018).
403    41.    A. Bateman, S. R. Eddy, C. Chothia, Members of the immunoglobulin superfamily in bacteria.
404           *Protein Sci* **5**, 1939-1941 (1996).
405    42.    L. Aravind, E. V. Koonin, Gleaning non-trivial structural, functional and evolutionary information
406           about proteins by iterative database searches. *J Mol Biol* **287**, 1023-1040 (1999).
407    43.    J. D. Mijnes *et al.*, Structure-function analysis of the gE-gI complex of feline herpesvirus:
408           mapping of gI domains required for gE-gI interaction, intracellular transport, and cell-to-cell
409           spread. *J Virol* **71**, 8397-8404 (1997).

410  44.  F. Deryckere, H. G. Burgert, Early region 3 of adenovirus type 19 (subgroup D) encodes an HLA-
411       binding protein distinct from that of subgroups B and C. *J Virol* **70**, 2832-2841 (1996).
412  45.  H. S. Ginsberg *et al.*, Role of early region 3 (E3) in pathogenesis of adenovirus disease. *Proc Natl*
413       *Acad Sci U S A* **86**, 3823-3827 (1989).
414  46.  R. H. Xu *et al.*, Antibody inhibition of a viral type 1 interferon decoy receptor cures a viral
415       disease by restoring interferon signaling in the liver. *PLoS Pathog* **8**, e1002475 (2012).
416  47.  J. K. Taylor *et al.*, Severe Acute Respiratory Syndrome Coronavirus ORF7a Inhibits Bone Marrow
417       Stromal Antigen 2 Virion Tethering through a Novel Mechanism of Glycosylation Interference. *J*
418       *Virol* **89**, 11820-11833 (2015).
419  48.  M. Oostra, C. A. de Haan, P. J. Rottier, The 29-nucleotide deletion present in human but not in
420       animal severe acute respiratory syndrome coronaviruses disrupts the functional expression of
421       open reading frame 8. *J Virol* **81**, 13876-13888 (2007).
422  49.  S. K. Lau *et al.*, Severe Acute Respiratory Syndrome (SARS) Coronavirus ORF8 Protein Is Acquired
423       from SARS-Related Coronavirus from Greater Horseshoe Bats through Recombination. *J Virol* **89**,
424       10532-10547 (2015).
425  50.  A. von Brunn *et al.*, Analysis of intraviral protein-protein interactions of the SARS coronavirus
426       ORFeome. *PLoS One* **2**, e459 (2007).
427
428

# SUPPLEMENTAL INFORMATION

**Supplementary Figure S1. Genome comparison analysis of SARS-related CoVs.**



Similarity Plot of SARS-related CoVs against human SARS-CoV-2 Wuhan-Hu-1 genome (NC_045512.2) based on a multiple sequence alignment of the whole genomes. Each point represents a different slicing window size from the alignment with a different step size between each point. For each plot, the window size and step size are shown in the top left. Horizontal bars above the top plot correspond to the different open reading frames of the SARS-CoV-2 genome (NC_045512.2). Each different colored line corresponds to the nucleotide similarity between the human SARS-CoV-2 genome (NC_045512.2) and the respective SARS-related CoV genome. The red arrows and solid lines surround regions which display recombination within the SARS-related CoV genomes. The single red arrows point to specific regions of recombination.

**Supplementary Figure S2. Full length multiple sequence alignment of ORF7a, ORF8-Ig and ORF7a/8-like proteins.**
Each sequence in the MSA was labelled by its species abbreviation followed by its isolation and NCBI accession number. The predicted secondary structure is shown above the alignment and the consensus is shown below the alignment, where h stands for hydrophobic residues, s for small residues, and p for polar residues. The characteristic signal peptide, TM region and a stretch of basic residues are also labeled.

**Supplementary Figure S3. Structural analysis of CoV Ig domains.**



ORF7a-Ig domain
(PDB: 1xak_A)

Model of ORF8-Ig domain

Location of highly variable
residues of ORF8-Ig domain
(sticks in green and red )

Surface view of ORF7a-Ig domain
(PDB: 1xak_A)

Surface view of ORF8-Ig domain

Location of hypervariable residue Leu84
on substrate-binding groove of the ORF8-ig domain
(stick in red)

## Supplementary Figure S4. Multiple sequence alignment of alphacoronavirus E3-CR1-like Ig domain proteins and their related Ig domains identified by profile-profile searches.

```
                             Signal peptide              β1            β2            β3                          β4            β5            β6
Secondary Structure
Bat CoV YN2012_Ra3376 QBP43259.1   - M L L L N V T T I V G C I I T A S L Q V D Q L T I T A N T L E R V D F H L A S S R K V C Q S S H S F N P T R Y K C N N N T L T L S V F E - - G Y K E T F E V R C F V K D - - T F Y R G
Bat CoV YN2012_Ra4259 QBP43281.1   - M L L L N V T T I I G C M L T A T I H S N Y L S I S A D T L E R V D F Y L A S G G K V C Q H S H S F N P T K Y K C E N N T L S V F V F A - - G Y K Q T F E V R C F V K D - - S F E S G
Bat CoV YN2012_Ra4125 QBP43270.1   - M L L L N V T T I I G C M L T A T I H S N Y L S I S A D T L E R V D F Y L A S G R K V C Q H S H S F N P T K Y K C E N N T L S V L C L L - - D T N K L L R F V V L L K T - - L L R V V
Bat CoV HKU32_TLC26A QCX35170.1     - M S V I A C L V L V T V - - - - - T V S H L N L F V D F N G R V D Y Y L G S G R K I C S F G H S F N P Q L Y N C T N T S L S L D I W Q - - G Y I G E F S V H C V G D - - - S I V S N
Bat CoV YN2012_Ra4259 QBP43282.1   - - M F F L L L I V V P V - - - - - - S C C N V S L V Q Y - N T T F Q Y T I D G E Y - - K K V E W K Y N N T H L I C S D G E V Y E Q F N T T V K C D N I S L V F D I A N V T L P N V T I
Bat CoV YN2012_Ra4125 QBP43271.1   - - M F L L L I V V P V - - - - - - - S C C N V S L V Q Y - N T T F Q Y T I D G E Y - - K K V E W K Y N N T H L I C S D G E V Y E Q F N T T V K C D N I S L V F D I A N V T L P N V T I
Bat CoV YN2012_Ra3376 QBP43260.1   - - M L F I I L V I V A V V - - - - - D C C E V S V V P H - N S T F Q Y T I D G D Y - - K H I N W K Y N D S L I I C S D G L V Y Q Q F N L T V K C D N L S L E F D T L Y V S I P S I T I
Bat CoV YN2012_Ra1359l QBP43292.1  - M L V N C V Y S L I T V V T - - - L N C T W F L L D A N - T V T I S L N V S G A - - - Q R V D W F K N I T K K L C S D G H S F Q G D - - V F I C N K T T L T T Q L K V G - - S N V S F
Bat CoV HKU32_TLC26A QCX35171.1     M L S M L T P L I V I A L S A - - - L H C C S C C T P Y - - I N G S T V S V E S N Y - - K S V E W K F N - N S Y I C S G G V Y D T F N A T F T C E N S T L T G N Y S G - - I D V L Q I
consensus/85%                       . . . h h . h . . h h s h . . . . . . p s s . h s h s . . . . p h p h . h s s . . . . p p . p h p h p . s . h . C p s s . h . . . h . . . . . h p p . . . . . . s h . . . . . s . . p .
```

```
                             β7                                                                                    TM region
- - - - - - - - - - - - - - - - - - - Y V A I S Y S - - - - - - - - - - - - - - - - - - - - - - - - - - - - - N I V R P N T P L I V I C I V I F C V L C Y S S Y - - - - - -
- - - - - - - - - - - - - - - - - - - Y V S I S E S - - - - - - - - - - - - - - - - - - - - - - - - - - - - - S K S K P N I P L L L V C S F V A I A L S C S F Y F - - - - -
- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - T C P F L S L L S L N L T Y H C F - - - - - -
E C K T K K - G E A S S V T F N N T I T V V T T H S P - - - - - - - - - - - - - - - - - T T N P S K S L I P S T K R H Y Y F L L L A F I P P A W F A V L I I Y Y V N P S G
E C K T K K - G E A S S V T F N N T I T V V T T Q S P - - - - - - - - - - - - - - - - - T T N P S K S L I P S T K R H Y Y F L L L A F I P P A W F A V L I I Y Y V N P S G
E C K T N N - G T H S V V T L N T I D S V K Y L S T P - - - - - - - - - - - - - - - - - L Q F P E Y S L I P S T K R L Y Y Y L P L S L I P P A W F T V P L I H Y V D T S -
N F H A T H - G G - - - - - - - T V H N G V F N V S V Y E S T A I T Q P P P L S S L L S P L H H Q A V E Y P S T K R L Y Y L L P L A F V I P A W L V V L F I H Y A D F S K
E C K M H N G S S L S A I V N F T T Q A P P T T V F V - - - - - - - - - - - - - - - - - T T Q T S R L I P S T K R T H Y S L F V A F T I V - P V V L V F I H Y A D F S -
. . . . . . . . . . . . . . . . . . . s . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . s . . p . . . . l . h . h . . . . . h . h . . . . a . . . . . . .
```

>PF02440.15 ; Adeno_E3_CR1 ; Adenovirus E3 region protein CR1
Probab=90.89  E-value=2.8  Score=29.56  Aligned_cols=38  Identities=18%  Similarity=0.352  Sum_probs=0.0  Template_Neff=6.800

```
Q QBP43282.1      24 QYTIDGE-YKKVEWK----YNNTHLICSDGEVYEQFNTTVKCDNISLVF    67 (114)
Q Consensus       24 s~svd~~~y~rVdw~-----n~s~kICs~ghsfn~f~~~~kC~N~TLt~   67 (114)
                      .++..+. +..|.|+    -..+.++|.....+        |+|++.+|++
T Consensus       17 NvTL~g~~~~~v~Wyr~~~~~~~~~LC~~~~~~~~~~~~~C~~~nLtL    59 (88)
T E3TMR2_ADE16/1   17 TCTLQGPQEGHVTWWRIYDNGGFARPCDQPGTK------FSCNGRDLTI    59 (88)
T ss_pred            cEEEeCCCCCCeEEEEeccCCCcccccccCCCe------eEeCCCceEE
```

>6JXR_n T-cell surface glycoprotein CD3 zeta; IMMUNE SYSTEM Homo sapiens
Probability: 62.4%, E-value: 84, Score: 22.12, Aligned cols: 126, Identities: 11%, Similarity: 0.052,

```
Q Q_9601142      19 VQYNTTFQYTIDGEYK---KVEWKYNNTHLICSDGEVYEQFN------TTVKCDNISLVFDIANVTLPNVTIECKTKKGE    89 (148)
Q Consensus      19 ~~~n~t~~~tv~g~yk----VeWk~N~s~~iCSdG~vyq~fn-------t~~Cdn~tL~~~~~~vs~p~vtieck~~~G~    89 (148)
                     ...+.+++-.+.+.+.    .+.|.++......-+.........    ..-..-..+|++.....+-..-.+.|.+.+..
T Consensus     135 ~g~~~~l~C~~~~~~~~~~~~~~W~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~l~i~~~~~d~~g~y~C~~~~~~   214 (291)
T 6JXR_n        135 HTQKATLVCLATGFYPDHVELSWWVNGKEVHSGVSTDPQPLKEQPALNDSRYCLSSRLRVSATFWQNPRNHFRCQVQFYG   214 (291)

Q Q_9601142      90 ASSVTFNNTITVVTTHSPTTNPSKSLIPSTKRH----------YYFLLLAFIPPAWFAVLIIYYV   144 (148)
Q Consensus      90 ~~~~~vnn~~~~~tt~~p~~~p~~sLiPSTKR~----------yY~L~lafi~paw~~V~iihYv   144 (148)
                     .................+.......--+.....             .+.+.++++.-..++++++-.+
T Consensus     215 ~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~i~~~~   279 (291)
T 6JXR_n  215       LSENDEWTQDRAKPVTQIVSAEAWGRADCGFTSESYQQGVLSATILYEILLGKATLYAVLVSALV   279 (291)
```
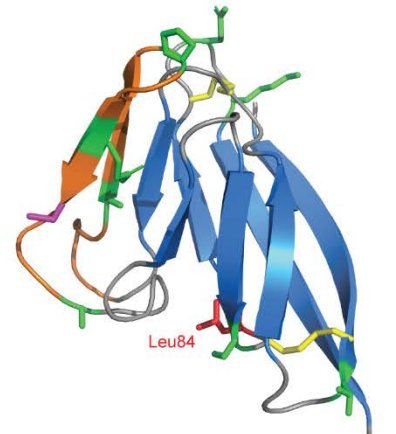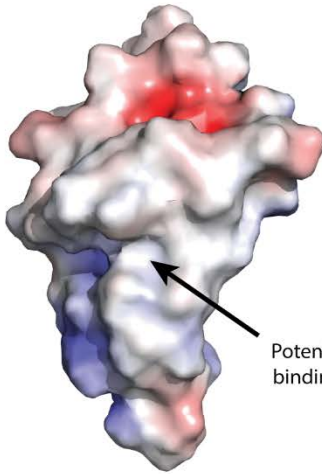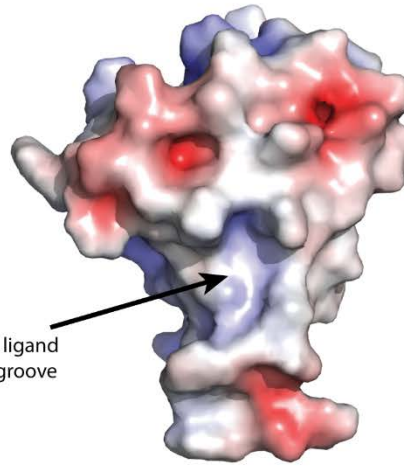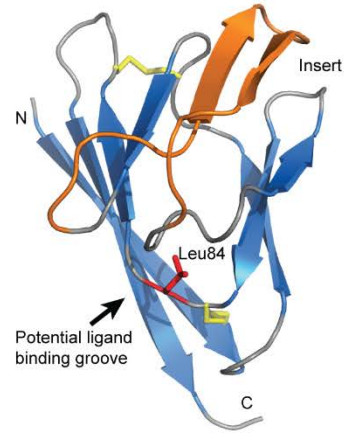
>3OQ3_B Interferon alpha-5, IFN-alpha/beta binding protein; Ectromelia, Mousepox Virus, Moscow strain; Mus musculus
Probability: 51.52%, E-value: 150, Score: 21.59, Aligned cols: 82, Identities: 12%, Similarity: 0.143,

```
Q Q_9601142      19 VQYNTTFQYTIDGEYK-----KVEWKYNNTHLICSDGEVYEQFNTTVKCDNISLVFDIANVTLPNVTIECKTKKGEASSV    93 (148)
Q Consensus      19 ~~~n~t~~~tv~g~yk------VeWk~N~s~~iCSdG~vyq~fn~t~~Cdn~tL~~~~~~vs~p~vtieck~~~G~~~~~    93 (148)
                     .....+++-.+.|...     .+.|.+++..+-.+...+.........-.+.+|++-.....--.-...|.+.|+.....
T Consensus     242 ~g~~~~l~C~~~~~~~~~~~~~~W~~~~~~~~~~~~~~~~~~~~~~~~~~l~i~~~~~~d~G~~Y~C~a~n~~g~~~   321 (329)
T 3OQ3_B        242 IGEPANITCTAVSTSLLVDDVLIDWENPSGWIIGLDFGVYSILTSSGGITEATLYFENVTEEYIGNTYTCRGHNYYFDKT   321 (329)

Q Q_9601142      94 TFNNTIT   100 (148)
Q Consensus      94 ~vnn~~~   100 (148)
                     .-.++..
T Consensus     322 ~~~~l~v   328 (329)
T 3OQ3_B        322 LTTTVVL   328 (329)
```

**Supplementary Table S1. Detailed information of Human SARS-CoV-2 Wuhan-Hu-1 genome and other SARS-related genomes that were used in this study (Figure 1 & Figure S1).**

| Organism | Host | NCBI ID | Year | Citation |
|---|---|---|---|---|
| Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) isolate Wuhan-Hu-1 | Human | NC_045512.2 | 2020 | A novel coronavirus associated with a respiratory disease in Wuhan of Hubei province, China (Unpublished) |
| Bat SARS-like coronavirus isolate bat-SL-CoVZC45 | Bat | MG772933.1 | 2018 | (1) |
| Bat SARS-like coronavirus isolate bat-SL-CoVZXC21 | Bat | MG772934.1 | 2018 | (1) |
| Bat coronavirus isolate RaTG13 | Bat | MN996532.1 | 2020 | Not Available |
| Severe acute respiratory syndrome-related coronavirus | Human | NC_004718.3 | 2003 | (2) |
| Severe acute respiratory syndrome-related coronavirus isolate F46 | Bat | KU973692.1 | 2017 | Identification of a new intermediate virus between bat-CoVs and SARS-CoVs from least horseshoe bats in China (Unpublished) |
| Bat SARS-like coronavirus YNLF_31C | Bat | KP886808.1 | 2015 | Not Available |
| Rhinolophus affinis coronavirus isolate LYRa11 | Bat | KF569996.1 | 2014 | (3) |
| Bat SARS coronavirus HKU3-7 | Bat | GQ153542.1 | 2010 | (4) |
| BtRs-BetaCoV/HuB2013 | Bat | KJ473814.1 | 2015 | (5) |
| Bat SARS-like coronavirus RsSHC014 | Bat | KC881005.1 | 2013 | (6) |
| Bat SARS-like coronavirus isolate Rs4231 | Bat | KY417146.1 | 2017 | (7) |

**Supplementary Table S2. Summary of representatives of viral Ig domain proteins which were identified in this study.**

| Family | Genus | Organism | NCBI ID | pfam ID | Domain Famliy | Presence of Signal Peptide (Y/N)* | Number of Ig-like domain | Number of TM region** | Distinct Relative PDB Structure# |
|---|---|---|---|---|---|---|---|---|---|
| Coronaviridae | Beta-coronavirus | SARS-CoV-2 | YP_009724395.1 | PF08779 | SARS_X4 | Y | 1 | 1 | 1XAK_A |
| | | | YP_009724396.1 | PF12093 | Corona_NS8 | Y | 1 | 0 | 1XAK_A |
| | | SARS-CoV | NP_828857.1 | PF08779 | SARS_X4 | Y | 1 | 1 | 1XAK_A |
| | | | NP_828876.1 NP_828877.1 | PF08779 | Corona_NS8 | Y | 1 | 0 | 1XAK_A |
| | Alpha-coronavirus | Bat coronavirus | QBP43259.1 | n/a | Adeno_E3_CR1-like | Y | 1 | 1 | 5XMZ_A |
| | | | QBP43265.1 | PF08779 | SARS_X4 | Y | 1 | 1 | 1XAK_A |
| Adenoviridae | Mast-adenovirus | Human adenovirus 7d | AAF14132.1 | PF02440 | Adeno_E3_CR1 | Y | 1 | 1 | 6JXR_d |
| | | Human adenovirus 23 | AFK92306.1 | PF02440 | Adeno_E3_CR1 | Y | 3 | 1 | 3J8F_7 |
| | | Human adenovirus 21 | AAW33363.1 | PF04881 | Adeno_GP19K | Y | 1 | 1 | 5IRO_P |
| Herpesviridae | Mardivirus | Gallid alphaherpesvirus 2 | YP_001034013.1 | PF02480 | Herpes_gE | Y | 1 | 1 | 2GJ7_F |
| | | | YP_001034012.1 | PF01688 | Herpes_gI | Y | 1 | 1 | 5OR7_C |
| | | | YP_001033973.1 | PF02124 | Marek_A | Y | 3 | 1 | 3J8F_7 |
| | Simplexvirus | Macacine alpha-herpesvirus 1 | NP_851925.1 | PF01537 | Herpes_glycop_D | Y | 1 | 1 | 4MYV_A |
| | Rhadinovirus | Human gamma-herpesvirus 8 | YP_001129350.1 | PF02960 | K1 | Y | 2 | 1 | 5D6D_C |
| | Cytomegalo-virus | Panine beta-herpesvirus 2 | NP_612760.1 | PF16758 | UL141 | Y | 1 | 1 | 4JM0_B |
| | | | NP_612778.1 | PF05963 | Cytomega_US3 | Y | 1 | 1 | 1IM3_P |
| | | Human beta-herpesvirus 5 | ABV71546.1 | PF17622 | UL16 | Y | 1 | 1 | 2WY3_B |
| | | Aotine beta-herpesvirus 1 | YP_004940175.1 | PF08001 | CMV_US | Y | 1 | 2 | 1IM3_P |
| Poxviridae | Orthopoxvirus | Variola virus | NP_042191.1 | PF08204 | V-set_CD47 | Y | 1 | 5 | 5OR7_C |
| | | Ectromelia virus | 3OQ3_B | PF13895 | ig | Y | 3 | 0 | 3OQ3_B |
| Phenuiviridae | Goukovirus | Cumuto virus | YP_009664616.1 | PF07245 | Phlebovirus_G2 | Y | 4 | 1 | 6F8P_A 6EGU_B |

* Signal Peptide Prediction was conducted by SignalP-5.0 program (*8*).

** Transmembrane (TM) region predictions were conducted by TMHMM Server (*9*).

# The PDB structures which display similarity with the respective viral Ig domains identified by pofile-profile comparisons (*10*).

**REFERENCES**

1.    D. Hu *et al.*, Genomic characterization and infectivity of a novel SARS-like coronavirus in Chinese bats. *Emerging microbes & infections* **7**, 1-10 (2018).
2.    M. A. Marra *et al.*, The genome sequence of the SARS-associated coronavirus. *Science* **300**, 1399-1404 (2003).
3.    B. He *et al.*, Identification of diverse alphacoronaviruses and genomic characterization of a novel severe acute respiratory syndrome-like coronavirus from bats in China. *Journal of virology* **88**, 7070-7082 (2014).
4.    S. K. Lau *et al.*, Ecoepidemiology and complete genome comparison of different strains of severe acute respiratory syndrome-related Rhinolophus bat coronavirus in China reveal bats as a reservoir for acute, self-limiting infection that allows recombination events. *Journal of virology* **84**, 2808-2819 (2010).
5.    Z. Wu *et al.*, Deciphering the bat virome catalog to better understand the ecological diversity of bat viruses and the bat origin of emerging infectious diseases. *The ISME journal* **10**, 609-620 (2016).
6.    X.-Y. Ge *et al.*, Isolation and characterization of a bat SARS-like coronavirus that uses the ACE2 receptor. *Nature* **503**, 535-538 (2013).
7.    B. Hu *et al.*, Discovery of a rich gene pool of bat SARS-related coronaviruses provides new insights into the origin of SARS coronavirus. *PLoS pathogens* **13**,  (2017).
8.    J. J. A. Armenteros *et al.*, SignalP 5.0 improves signal peptide predictions using deep neural networks. *Nature biotechnology* **37**, 420-423 (2019).
9.    A. Krogh, B. Larsson, G. Von Heijne, E. L. Sonnhammer, Predicting transmembrane protein topology with a hidden Markov model: application to complete genomes. *Journal of molecular biology* **305**, 567-580 (2001).
10.   J. Söding, A. Biegert, A. N. Lupas, The HHpred interactive server for protein homology detection and structure prediction. *Nucleic acids research* **33**, W244-W248 (2005).