

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16

Re-annotation of SARS-CoV-2 proteins using an HHpred-based approach opens new opportunities for a better understanding of this virus

Pierre Brézellec*^{1,2}

¹ Institut Systématique Evolution Biodiversité (ISYEB UMR 7205), Sorbonne Université, MNHN, CNRS, EPHE, UA, Paris, France.

² Université de Versailles Saint Quentin, 45 avenue des Etats Unis, 78000 Versailles, France.

*Corresponding author

Correspondence: pierre.brezellec@uvsq.fr

ABSTRACT

Since the publication of the genome of SARS-CoV-2 – the causative agent of COVID-19 – in January 2020, many bioinformatic tools have been applied to annotate its proteins. Although efficient methods have been used, such as the identification of protein domains stored in Pfam, most of the proteins of this virus have no detectable homologous protein domains outside the viral taxa. As it is now well established that some viral proteins share similarities with proteins of their hosts, we decided to explore the hypothesis that this lack of homologies could be, at least in part, the result of the documented loss of sensitivity of Pfam Hidden Markov Models (HMMs) when searching for domains in "divergent organisms". In order to improve the annotation of SARS-CoV-2 proteins, we used the HHpred protein annotation tool. To avoid "false positive predictions" as much as possible, we designed a robustness procedure to evaluate the HHpred results. In total, 6 robust similarities involving 6 distinct SARS-CoV-2 proteins were detected. Of these 6 similarities, 3 are already known and well documented, and

30 one is in agreement with recent crystallographic results. We then examined carefully the two
31 similarities that have not yet been reported in the literature. We first show that the C-terminal
32 part of Spike S (the protein that binds the virion to the cell membrane by interacting with the
33 host receptor, triggering infection) has similarities with the human prominin-1/CD133; after
34 reviewing what is known about prominin-1/CD133, we suggest that the C-terminal part of Spike
35 S could both improve the docking of Spike S to ACE2 (the main cell entry receptor for SARS-
36 CoV-2) and be involved in the delivery of virions to regions where ACE2 is located in cells.
37 Secondly, we show that the SARS-CoV-2 ORF3a protein shares similarities with human G
38 protein-coupled receptors (GPCRs) belonging mainly to the "Rhodopsin family"; ~~on the basis of
39 the literature, we then show that specific G protein-coupled receptors (GPCRs) of this family
40 are known to form ion channels; we emphasize this is consistent with a recent Cryo-EM
41 structure of SARS-CoV-2 ORF3a suggesting that it can form a non-selective Ca²⁺-permeable
42 cation channel; furthermore, we highlight that some of the GPCRs identified as sharing
43 similarities with ORF3a are targeted by antibodies in patients with COVID-19 and Long COVID;
44 suggesting that these similarities may trigger some of the observed autoimmune responses.~~ We
45 conclude that the approach described here (or similar approaches) opens up new avenues of
46 research to better understand SARS-CoV-2 and could be used to complement virus annotations,
47 particularly for less-studied viruses.

48 **Keywords:** Pfam Domains, HHpred, Hidden Markov Models (HMMs), Bioinformatics, Protein annotation,
49 SARS-CoV-2.

50

51

52

54 A significant fraction of the proteins expressed by viruses often lack homologs. These proteins are termed
55 "orphan" to emphasise that no homologs are detected, or "taxonomically restricted" to indicate that they
56 have no detectable homologs outside a given taxon (Kuchibhatla *et al.*, 2014). SARS-CoV-2 (Severe Acute
57 Respiratory Syndrome Coronavirus 2), the causative agent of COVID-19, is no exception. According to UniProt
58 (UniProt Consortium, 2021), this virus expresses 17 proteins (see Supplemental file 1 for more details). If we
59 consider the Pfam annotations (Mistry *et al.*, 2021, <http://pfam-legacy.xfam.org/>) of the proteins expressed
60 by this virus, we observe that *i/* 4 of these 17 proteins are not Pfam annotated, *ii/* the other 13 proteins are
61 annotated by a set of 40 domains, 39 of which are strictly associated with viruses (the Macro domain being
62 an exception to the rule). This clearly shows that SARS-CoV-2 domains are mostly similar to viral domains
63 (97.5% $((39/40) * 100)$ which are generally poorly annotated.

64 These results can be interpreted in two different (but complementary) ways:

65 1./ This virus, like many viruses, essentially contains virus-like proteins that are only present in viruses
66 and not elsewhere,

67 2./ As it has been established that *i/* some viral proteins show similarities to some proteins of their host
68 and that *ii/* this "molecular mimicry" is increasingly recognised (Elde & Malik, 2009), this lack of homologies
69 outside of viral taxa can also be seen, at least in part, as a consequence of weaknesses in annotation
70 methods.

71 It has been shown that HMMs stored in Pfam can lack sensitivity when searching for domains in
72 "divergent organisms" (where the relevant signals become too weak to be identified (Terrapon *et al.*, 2012)).
73 We thus decided here to explore the second way. We naturally turned to HHpred which is known to be an
74 efficient tool for remote protein homology detection and can be easily used via a fast server (Gabler *et al.*,
75 2020). HHpred offers many possibilities such as searching for homologs among all proteins in an organism.
76 HHpred is based on HHsearch and HHblits, which perform pairwise comparison of HMM profiles. ~~Given their
77 proven efficiency, HHsearch and HHblits have been used for some years to annotate viruses, and in particular
78 Coronaviruses (Forni *et al.*, 2022). They have obviously also been used to annotate proteins expressed by
79 SARS-CoV-2 (O'Donoghue *et al.*, 2021). However, the two previous works limited the homology search to
80 viral proteins. For our part, we focused on searching for homologs in human. Given their proven efficiency,
81 HHsearch and HHblits have been used for some years to annotate viruses, and in particular accessory
82 proteins of coronaviruses (Forni *et al.*, 2022). They have also been used to model proteins structures
83 expressed by SARS-CoV-2 (O'Donoghue *et al.*, 2021) using related 3D structures in the PDB, *i.e.*, structures
84 determined for other coronaviruses, such as SARS-CoV or MERS-CoV, as well as many structures from more
85 distantly related viruses, such as those causing polio or foot-and-mouth disease. However, the two previous
86 works limited the homology search to viral proteins. Here, using an available database of HMMs specific to
87 *Homo sapiens* proteins, we directly searched - using HHpred - for homologs of SARS-CoV-2 proteins in
88 human. Thus, what was previously achievable at the Pfam domain level (for instance) now extends to human
89 proteins.~~

90 To avoid "false positive predictions" as much as possible, we designed a procedure, mainly based on two
91 ideas suggested in (Gabler *et al.*, 2020) but not implemented, to assess the robustness of HHpred results.
92 Using HHpred and this procedure, we detected 6 robust similarities.

94 SARS-CoV-2 protein sequences

95 The 17 proteins studied in this article were extracted from UniProt (<https://www.uniprot.org/>, UniProt
 96 Consortium, 2021). UniProt provides polyproteins 1a (pp1a) and 1ab (pp1ab) as two separate entries. The
 97 pp1ab polyprotein is cleaved to form 15 shorter proteins; the first 10 proteins, *i.e.*, NSPs 1-10, are also
 98 cleaved from pp1a; NSPs 12-16 are unique to pp1ab. The list of proteins is given below. For each protein, we
 99 give its "Recommended Name", its "Short Name", its "AC - Uniprot ID", and its length:
 100 Replicase polyprotein 1a / pp1a / P0DTC1 - R1A_SARS2 / Length 4,405
 101 Replicase polyprotein 1ab / pp1ab / P0DTD1 - R1AB_SARS2 / Length 7,096
 102 Envelope small membrane protein / E; sM protein / P0DTC4 - VEMP_SARS2 / Length 75
 103 Membrane protein / M / P0DTC5 - VME1_SARS2 / Length 222
 104 Nucleoprotein / N / P0DTC9 - NCAP_SARS2 / Length 419
 105 Spike glycoprotein/ S glycoprotein / P0DTC2 - SPIKE_SARS2 / Length 1,273
 106 ORF3a protein/ ORF3a / P0DTC3 - AP3A_SARS2 / Length 275
 107 ORF3c protein / ORF3c / P0DTG1 - ORF3C_SARS2 / Length 41
 108 ORF6 protein / ORF6 / P0DTC6 - NS6_SARS2 / Length 61
 109 ORF7a protein / ORF7a / P0DTC7 - NS7A_SARS2 / Length 121
 110 ORF7b protein / ORF7b / P0DTD8 - NS7B_SARS2 / Length 43
 111 ORF8 protein / ORF8 / P0DTC8 - NS8_SARS2 / Length 121
 112 ORF9b protein / ORF9b / P0DTD2 - ORF9B_SARS2 / Length 97
 113 Putative ORF3b protein/ ORF3b / P0DTF1 - ORF3B_SARS2 / Length 22
 114 Putative ORF3d protein/ _ / P0DTG0 - ORF3D_SARS2 / Length 57
 115 Putative ORF9c protein / ORF9c / P0DTD3 - ORF9C_SARS2 / Length 73
 116 Putative ORF10 protein / ORF10 / A0A663DJA2 - ORF10_SARS2 / Length 38

117 Sequence similarity searches

118 For remote homology detection, we used HHpred (Gabler *et al.*, 2020). First, starting from single
 119 sequences or multiple sequence alignments (MSAs), it transforms them into a query HMM; using this HMM,
 120 it then searches the Uniclust database³⁰ and adds significantly similar sequences found to the query MSA for
 121 the next search iteration. This strategy is very effective in detecting remotely homologous sequences but, as
 122 the user guide points out (<https://github.com/soedinglab/hh-suite/wiki>), "the higher the number of search
 123 iterations, the greater the risk of non-homologous sequences or sequence segments entering the MSA and
 124 recruiting other sequences of the same type in subsequent iterations". To avoid this problem, we set the
 125 number of iterations to 0, *i.e.* the parameter "MSA generation iterations" was set to 0. The default settings
 126 were used for the other parameters. Note that we also briefly present in the Results section the HHpred
 127 results obtained using the default setting for "MSA generation iterations", *i.e.* 3 (iterations).

128 Finally, it is important to note here that we use HHpred to look for similarities independently of the
 129 mechanisms underlying these similarities, *i.e.* homologies, horizontal transfers (e.g. obtained by
 130 "recombination" between SARS-CoV-2 and its current host, between ancestors of SARS-CoV-2 and their
 131 hosts, between SARS-CoV-2 and another virus, *etc.*), convergent evolutions, *etc.*

132 Procedure for assessing the robustness of HHpred results

133 According to (Gabler *et al.*, 2020), when the reported probability value for a hit is greater than 95%,
 134 homology is highly probable. Since viral and human proteins are being compared here, it can be assumed
 135 that the 95% threshold is too high to detect similarities. In order to be more sensitive, while controlling
 136 specificity (*i.e.* avoiding "false positive predictions"), we have devised a procedure that we describe below. Its

137 purpose is to assess the robustness of the results provided by HHpred. It is based on two ideas described in
138 (Gabler *et al.*, 2020) (section "Understanding Results") but not taken into account in HHpred.

139 This procedure is divided into 4 steps. From an algorithmic point of view, this procedure can be described
140 as a "greedy search algorithm". It is performed for each protein expressed by the SARS-CoV-2 virus (see
141 Figure) (note that this was done by hand, as there were few data to process):

142 1./ For a given SARS-CoV-2 protein, hereafter referred to as "query", HHpred is run (using the default
143 parameters, except for the "MSA generation iterations" parameter which we set to 0, see section above) on
144 the *Homo sapiens* proteome of HHpred.

145 2.1/ The examination of the results provided by HHpred starts with the probability threshold of 0.95. Hits
146 with a probability greater than or equal to 0.95 are selected. If no hits meet this constraint, the threshold is
147 successively lowered to 0.9, 0.85 and finally to 0.80. As soon as a threshold satisfies the constraint (*i.e.* there
148 is at least one hit with a probability greater than or equal to the threshold), all hits above the threshold are
149 selected. If no threshold satisfies the constraint, we consider that no similarity between the query and the
150 human proteins can be detected.

151 2.2./ All previously selected hits are collected in a list and ranked from highest to lowest probability. The
152 best hit is then used as a seed to build a family of hits as follows: hits located at the same position as this best
153 hit on the query sequence and of similar size to it (+5 amino acids for a best hit of length < 150, and +-15 for
154 a best hit of length > 150) feed the family under construction and are removed from the list ; hits that
155 overlap the seed are also removed from the list. The highest hit in the updated list is used as the "new" seed
156 and the process continues until the list is empty. As it is possible for a protein to have only one homolog in
157 human, families of singletons are not excluded.

158 3.1/ The query is then run on four HHpred proteomes, called "test" proteomes, corresponding to the
159 following four species: *Arabidopsis thaliana*, *Drosophila melanogaster*, *Escherichia coli* and *Haloferax volcanii*
160 (an archaea).

161 3.2/ For each species, the following step is performed:

162 First, the hits whose probability is greater than or equal to the previously selected threshold (see 2.1) are
163 selected. Then, depending on their size and location on the query sequence, they are assigned, if possible, to
164 a previously built family (see 2.2).

165 At the end of step 3, a family is thus made up of hits belonging at least to *Homo sapiens* and possibly to
166 *Arabidopsis thaliana*, *Drosophila melanogaster*, *Escherichia coli* or *Haloferax volcanii*. If a family includes only
167 human proteins, the robustness assumption can neither be rejected nor established. In this case, the
168 threshold is lowered and step 2 is performed again.

169 4./ For each family, InterPro annotations (Blum *et al.*, 2020) of proteins associated with hits are collected
170 and inspected manually (in particular the part of these proteins that corresponds to the hits). If the
171 annotations of the human proteins are similar to each other and to all proteins from at least one other
172 organism, this family/similarity is considered "robust"; these annotations are then associated with the
173 corresponding part of the viral protein (the query) ; if not, no similarities can be identified, and we consider
174 that no similarity between the query and the human proteins can be detected.

175 It should be noted that when the threshold of 0.8 is reached and it is not possible to reject or establish
176 the robustness hypothesis, an in-depth examination of the results is carried out by relaxing the constraints *i*/
177 on the probability threshold, which is then set to 0.5 (in accordance with the HHpred documentation which
178 states that "typically, a match should be seriously considered if it has a probability value >50%"), and *ii*/
179 on the size and location of hits ; the annotations of the proteins found by relaxing the constraints are then
180 examined; if at least 90% of human proteins are similarly annotated and these are also similarly annotated to
181 100% of the proteins of at least one other organism, this family/similarity is considered "robust"; these
182 annotations are then associated with the corresponding part of the viral protein (the query).

183 The similarities identified at the 0.95 and 0.9 probability levels will be labeled by "highly robust"; the
184 similarities identified at the 0.85 and 0.8 probability levels will be labeled by "very robust"; finally, the
185 similarities identified during the relaxation stage of constraints will be labeled by "quite robust".

186 Note: only proteins beginning with the prefix NP are considered in the analysis. XP records (proteins) are
187 not curated and are therefore not considered here; furthermore, proteins identified by HHpred that do not
188 have a match in "UniProtKB reviewed (Swiss-Prot)" (name and size in amino acids) were not considered
189 either.

190 Results

191 ~~In the following sections, we present the main results of our study, i.e., the list of the 6 robust similarities~~
192 ~~we have identified. In our study, we identified a list of 6 robust similarities. We focus here on the two~~
193 ~~similarities not yet documented in the literature.~~ For reasons of clarity, for each family/similarity considered
194 here, only the best hit in each organism is provided. All results can be found in Supplemental file 2 (this file
195 contains a condensed version of the results produced by HHpred which are enriched by the InterPro
196 annotations). The raw HHpred results are stored in a separate gzip file called Supplemental file 4.

197 Note that the Pfam annotations of the proteins come from the InterPro or Pfam legacy ([http://pfam-
198 legacy.xfam.org/](http://pfam-legacy.xfam.org/)) websites; the two sites generally give similar predictions; however, the domain boundaries
199 may sometimes differ very slightly.

200 ~~NSP2 harbors a "Casein kinase II regulatory subunit" domain (very robust similarity)~~

201 ~~NSP2 is derived from polyprotein 1a (181-818). The length of this protein is 638 A.A.~~

202 ~~At the 0.85 probability threshold, only one human protein shares similarity with NSP2. Specifically, the~~
203 ~~151-195 part of NSP2 is similar to the 101-142 part of the human protein "Casein kinase II subunit beta"~~
204 ~~(CSK2B_HUMAN/NP_001311, length = 215 A.A.). The 109-140 part of the human "Casein kinase II subunit~~
205 ~~beta" protein is annotated with the PROSITE "Casein kinase II regulatory subunit signature" motif; in~~
206 ~~addition, parts 105-126 and 127-148 of this protein are annotated with the PRINTS motif "CASNKINASEII".~~
207 ~~This suggests that the 151-195 part of NSP2 could also be a "casein kinase II regulatory subunit signature".~~
208 ~~Note that when the "MSA generation iterations" parameter is set to 3 (default setting), no significant results~~
209 ~~are obtained (the probability of the best hit is 0.32).~~

210 ~~For the given threshold of 0.85, this part of NSP2 is also similar to a part of an *Arabidopsis thaliana*~~
211 ~~protein. Specifically, the 151-194 part of NSP2 is similar to the 182-222 part of the *Arabidopsis thaliana*~~
212 ~~protein "Casein kinase II subunit beta" (CSK2D_ARATH/NP_191584.1, length = 276 AA). The 190-221 part of~~
213 ~~the latter is annotated with the PROSITE motif "Casein kinase II regulatory subunit signature"; in addition;~~
214 ~~parts 186-207 and 208-229 of this protein are annotated with the PRINTS motif "CASNKINASEII".~~

215 ~~This strongly suggests that NSP2 carries a "regulatory subunit signature of casein kinase II".~~

216 ~~NSP3 harbors a Macro domain (highly robust similarity)~~

217 ~~NSP3 is derived from polyprotein 1a (819-2763). The length of this protein is 1945 A.A.~~

218 ~~For the 0.95 probability threshold, 7 human proteins share similarity with NSP3. The best match is the~~
219 ~~human "Core histone macro H2A.2" protein (H2AW_HUMAN/NP_061119, length = 372) whose 187-371 part~~
220 ~~is similar to the 210-377 part of NSP3, i.e. the 1029-1197 part of polyprotein 1a. The 187-371 region of this~~
221 ~~human protein contains the Pfam Macro domain (216-329). This suggests that the 210-377 part of NSP13~~
222 ~~also shares similarity with the Macro domain. Note that when the "MSA generation iterations" parameter is~~
223 ~~set to 3 (default setting), similar results are obtained.~~

224 ~~For the probability threshold considered (i.e., 0.95), one *Escherichia coli* protein shares similarity with~~
225 ~~NSP3: the "O-acetyl-ADP-ribose deacetylase" protein (YMDB_ECOLI/NP_415563, length = 177) whose 3-166~~

226 part is similar to the 218-367 part of NSP3; the 218-367 region of this bacterial protein contains the Pfam
227 Macro domain (21-137).

228 The above strongly suggests that NSP3 hosts a Macro domain.

229 **NSP13 harbors AAA domains (highly robust similarity)**

230 NSP13 is derived from polyprotein 1ab (5325-5925). The length of this protein is 601 A.A.

231 At the 0.95 probability level, many human proteins share similarity with NSP13. The best match is the
232 human "DNA-binding protein SMUBP-2" (SMBP2_HUMAN/NP_002171, length = 993) whose 207-618 part is
233 similar to the 275-582 part of NSP13. The 207-618 part of this human protein is involved in two Pfam
234 domains, namely AAA_11/191-411 and AAA_12/418-615, which are both members of the P-loop NTPase clan
235 (CL0023). This suggests that the 275-582 part of NSP13, i.e. the 5600-5907 part of polyprotein 1ab, harbours
236 AAA domains. Note that when the "MSA generation iterations" parameter is set to 3 (default setting), similar
237 results are obtained.

238 The previously considered part of NSP13 is similar to three *Arabidopsis thaliana* proteins. The best match
239 is the *Arabidopsis thaliana* "probable helicase" protein (MAA3_ARATH/NP_001329005, length = 818) whose
240 273-734 part is similar to the 275-581 part of NSP13. The 273-734 part of this plant protein is involved in
241 three Pfam domains, namely AAA_11/257-436 + AAA_11/451-526 + AAA_12/534-731, which are members
242 of the P-loop NTPase clan (CL0023). This result is in agreement with what has been found in human.

243 Our results strongly suggest that the 275-582 part of NSP13 hosts AAA domains.

244 **NSP16 is a methyltransferase (highly robust similarity)**

245 NSP16 is derived from polyprotein 1ab (6799-7096). The length of this protein is 298 A.A.

246 At the 0.95 probability level, two human proteins share similarity with NSP13. The best match is the "pre-
247 rRNA 2'-O-ribose RNA methyltransferase FTSJ3" protein (SPB1_HUMAN/NP_060117, length = 847). Its 31-217
248 part is similar to the 46-230 part of NSP16, i.e. the 6845-7029 part of polyprotein 1ab. The 31-217 part of this
249 human protein corresponds quite well to the Pfam "FtsJ like methyltransferase" domain, FtsJ/21-207. Note
250 that when the "MSA generation iterations" parameter is set to 3 (default setting), similar results are
251 obtained.

252 The part of NSP16 considered above is similar to two *Drosophila melanogaster* proteins. The best match
253 is the fly protein "Putative tRNA (cytidine(32)/guanosine(34) 2'-O) methyltransferase 1"
254 (TRM71_DROME/NP_650590, length = 302) whose 28-211 part is similar to the 46-215 part of NSP16. The
255 28-211 part of this fly protein corresponds quite well to the Pfam "FtsJ like methyltransferase" domain,
256 FtsJ/21-207. This result is in agreement with what was found in human.

257 This strongly suggests that NSP16 is a methyltransferase.

258 **Spike S harbors a part of a "Prominin domain" (highly robust similarity)**

259 The length of this protein is 1273 A.A.

260 At the 0.90 probability level, 2 human proteins share similarity with Spike S (prominin-1 and prominin-2
261 proteins). The best match is human prominin-1 (PROM1_HUMAN/NP_006008, length = 865). Its 186-482 part
262 is similar to the 908-1254 part of Spike S; the 186-482 part of this human protein is included in the Pfam
263 "Prominin" domain, Prominin/19-820. Note that when the "MSA generation iterations" parameter is set to 3
264 (default setting), similar results are obtained.

265 For the given threshold of 0.90, one fly protein annotated with the Prominin domain of Pfam shares
266 similarities with Spike S: the fly protein "Prominin-like protein" (PROML_DROME/NP_001261351.1, length =
267 1013) whose 235-534 part is similar to the 911-1254 part of Spike S; the 235-534 part of this fly protein is
268 included in the "Prominin" domain of Pfam, Prominin/76-881.

269 This strongly suggests that Spike S hosts part of the "Prominin domain".

270 ORF3a has similarities with some "G Protein-Coupled Receptors" (quite robust similarity)

271 The length of this protein is 275 A.A.

272 At the 0.80 probability level, a human protein shares similarity with ORF3a, the human "lutein-
273 choriogonadotropic hormone receptor" (LSHR_HUMAN/NP_000224, length = 699). Its 537-693 part is similar
274 to the 41-183 part of ORF3a. A large part of this 537-693 region is included in the "7 transmembrane
275 receptor (rhodopsin family)" Pfam domain, *i.e.* 7tm_1/376-623. Note that when the "MSA generation
276 iterations" parameter is set to 3 (default setting), no significant results are obtained (the probability of the
277 best hit is 0.66).

278 For the given threshold of 0.80, no similarity is detected with proteins belonging to the 4 "test"
279 proteomes. However, a number of factors support this similarity when certain constraints are relaxed (see
280 Materials and methods):

281 Looking at the list of hits found by HHpred between ORF3a and the human proteome (see Supplemental
282 file 2), it is immediately obvious that the vast majority of human proteins found are G Protein-Coupled
283 Receptors (GPCRs). Indeed, it appears that out of 28 hits, 26 concern GPCRs (26/28 = 0.928), while the other
284 two correspond to transmembrane segments of proteins that are not linked to GPCRs.

285 Considering the fly proteome and applying the same methodology as previously used in human, it
286 appears that out of 3 hits, 3 concern GPCRs (see Supplemental file 2).

287 Overall (see Materials and Methods), this suggests that the similarity found is quite robust and that
288 ORF3a shares similarities with human GPCRs.

289

Discussion

290 The documented loss of sensitivity of Pfam HMMs when searching for domains in "divergent organisms"
291 (Terrapon *et al.*, 2012) prompted us to use HHpred (Gabler *et al.*, 2020) to annotate SARS-CoV-2 proteins.
292 Given a query sequence, this annotation tool offers the possibility to search for homologs among all proteins
293 in an organism. Each protein in the organism is represented by an HMM built according to a different
294 strategy than that used by Pfam (for more details, see the section "Creating custom databases" in the user
295 guide (<https://github.com/soedinglab/hh-suite/wiki>)). We speculated that this difference might give HHpred
296 the ability to discover similarities not detectable by Pfam (it should be noted that a theoretical comparison
297 between the Pfam and HHpred HMMs, as well as a full empirical comparison, is beyond the scope of this
298 paper).

299 To avoid as much as possible false predictions when using HHpred, we decided to disable its first step
300 which is based on an iterative search strategy. Indeed, the greater the number of search iterations, the
301 greater the risk of recruiting non-homologous sequences in the following iterations (see Materials and
302 Methods). Furthermore, in addition to the probability assigned by HHpred to each hit, we decided to
303 evaluate the robustness of these latter. Our evaluation procedure is based on two unimplemented ideas
304 described in (Gabler *et al.*, 2020) and can be summarized as follows (see Materials and Methods for more
305 details ; see also Figure):

306 A probability threshold is set; the starting value is 0.95 (according to (Gabler *et al.*, 2020), when the
307 probability of a hit is greater than 95%, homology is highly probable). Each viral protein ("query" sequence) is
308 compared to the human proteome using HHpred; all hits with a probability above the chosen threshold are
309 selected (if no hit meets this criterion, the threshold is successively lowered to 0.9, 0.85 and 0.80) ; if all hits
310 of similar size located at the same position on the query sequence (*i.e.*, a family of homologous hits) are
311 annotated with the same InterPro domain (Blum *et al.*, 2020), their probability of actually being homologous
312 to the query is very high ("Check relationships among top hits", first idea from (Gabler *et al.*, 2020)); the
313 query is then run on a set of "test" proteomes to check whether similarly annotated homologous hits are
314 returned ("Check if you can reproduce the results with other parameters", second idea of (Gabler *et al.*,

2020)); if so, a family of homologous hits defined a "robust similarity"; if not, we consider that no similarities can be identified. Note that when a family includes only human proteins, the robustness assumption can neither be rejected nor established; in this case, the threshold is lowered and the study is carried out again. It should be also noted that when the threshold of 0.8 is reached and it is not possible to reject or establish the robustness hypothesis, a thorough examination of the results is carried out by relaxing the constraints (mainly on the size, location and/or probability associated with the hits, see Materials and Methods for more details). Similarities identified at the 0.95 and 0.9 probability levels are labeled "highly robust"; similarities identified at the 0.85 and 0.8 probability levels are labeled "very robust"; finally, the similarities identified when certain constraints are relaxed are described as "quite robust".

The organisms used to evaluate the HHpred results are *Arabidopsis thaliana*, *Drosophila melanogaster*, *Escherichia coli* and *Haloferax volcanii* (an archaea). Note that, in order to potentially increase the identified similarities, we would have liked to include proteomes from organisms closer to humans in our study. Unfortunately, the online server currently does not offer the option to use such proteomes. To successfully accomplish this task, it is necessary to perform the local installation of the free HH-suite software and build these proteomes using this software. This work needs to be done (future works).

Below we present a summary of our results.

We subjected the 17 proteins of the SARS-CoV-2 proteome (see Materials & Methods and Results sections) to our annotation procedure. UniProt considers polyproteins 1a (pp1a) and 1ab (pp1ab) as two separate entries; polyprotein pp1ab is proteolytically cleaved to form 15 shorter proteins; the first 10 proteins (NSP1, ..., NSP10) are also cleaved from pp1a; NSP12, ..., NSP16 are unique to pp1ab. We therefore subjected 30 proteins to our evaluation procedure.

No "robust" similarities were found for the following 24 proteins

NSP1, NSP4-10, NSP12, NSP14-15, Nucleoprotein, Envelope small membrane, Membrane Protein M, ORF3B, ORF3C, ORF3D, ORF6, ORF7a, ORF7b, ORF8, ORF9b, ORF9C, ORF10.

A "highly robust" or "very robust" similarity, already documented in literature, was detected on the following 4 proteins

~~In a more interesting manner, we have shown that part 151-195 of NSP2, i.e. part 332-376 of polyprotein 1a, contains a "signature of the beta subunit of casein kinase II".~~

~~NSP3 harbors a Macro domain; NSP13 harbors AAA domains; NSP16 is a methyltransferase. As these similarities are well documented and widely discussed, the interested reader is invited to consult the InterPro annotations.~~

~~NSP3 is a papain-like protease; we showed it harbors a Macro domain. NSP13 is a helicase; we provide evidence suggesting that it harbors AAA domains. NSP16 is a methyltransferase; we confirm that it harbors a "FtsJ-like methyltransferase" domain. As these similarities are well documented, the interested reader is invited to consult the InterPro annotations.~~

~~NSP2 is involved in the inhibition of the antiviral response and facilitates SARS-CoV-2 replication. We showed that part 151-195 of NSP2, i.e. part 332-376 of polyprotein 1a, contains a "signature of the beta subunit of casein kinase II".~~ According to PROSITE, such a domain could be involved in the binding of a metal such as zinc. Interestingly, the structure of the N-terminal part of NSP2 was recently solved (Ma *et al.*, 2021). It shows that NSP2 has three zinc fingers: Zn1, Zn2 and Zn3. Two Zn2 (resp. Zn3) binding sites are located at positions 161 and 164 (resp. at positions 190 and 193). Our prediction is therefore in agreement with this structure of the N-terminal domain of SARS-CoV-2 NSP2.

357 **A previously unknown "highly robust" similarity was detected on Spike S protein**

358 The Spike S protein (1273 A.A.) is composed of two subunits: the S1 subunit (14-685 residues), and the S2
359 subunit (686-1273 residues), which are responsible for receptor binding and membrane fusion respectively
360 (Huang *et al.*, 2020). We have shown that the 908-1254 part of the Spike S protein is similar to the 186-482
361 part of human prominin-1 (length = 865). This similarity encompasses the heptapeptide repeat 1 sequence,
362 *i.e.* HR1 (912-984 residues), HR2 (1163-1213 residues), the TM domain (1213-1237) and part of the
363 cytoplasmic domain (1237-1273) of the S2 subunit; however, it excludes the fusion peptide (FP) (788-806) of
364 S2 which plays an essential role in mediating membrane fusion. HR1 and HR2, which are part of the
365 similarity, have been shown to form a six-helix bundle that is essential for the fusion and viral entry function
366 of the S2 subunit (Xia *et al.*, 2020).

367 Recently, in searching for proteins involved in SARS-CoV-2 entry into host cells, (Kotani *et al.*, 2022) found
368 that the glycoprotein CD133, the other name for prominin-1, colocalises with ACE2 – the main cell entry
369 receptor for SARS-CoV-2 – bound to the Spike S protein in Caco-2 cells. They demonstrated that the SARS-
370 CoV-2 Spike protein exhibited increased binding capacity in cells co-expressing ACE2 and CD133, compared
371 to cells expressing ACE2 alone. In addition, they experimentally infected HEK293T cells with a SARS-CoV-2
372 pseudovirus and showed that infectivity was twice as high in HEK293T cells co-expressing CD133-ACE2 than
373 in HEK293T cells expressing ACE2 alone. They concluded that CD133, although not a primary receptor for the
374 SARS-CoV-2 Spike protein, is a cofactor (a co-receptor) that partially contributes to infection in the expressing
375 cells. All these results suggest that the C-terminal part of Spike S, which has similarities with prominin-1, may
376 be involved in the docking of Spike S to ACE2 (insofar as CD133 enhances the ability of Spike S to bind to
377 ACE2). This obviously remains to be demonstrated but is clearly an interesting avenue of research.

378 While considerable work has been done to characterise the cellular receptors and pathways mediating
379 virus internalisation, little is known about the onset of the infection process, which begins when the virus
380 comes into contact with the host cell surface; some studies have shown that viruses "diffuse" onto the
381 surface of host cells after "landing" on them; this process ranges from a random walk to a constrained
382 diffusion where the virus particles appear to be confined to a specific microdomain of the cell membrane
383 (Boulant *et al.*, 2015). From this point of view, it is interesting to note that it was recently shown by (Rouaud
384 *et al.*, 2022) that *i/* ACE2 concentrates at epithelial apical cell junctions in cultured epithelial cell lines, and
385 that *ii/* (Pinto *et al.*, 2022) showed that ACE2 and TMPRSS2 (which is used by SARS-CoV-2 for Spike S-protein
386 priming (Hoffmann *et al.*, 2020)) were localised at the plasma membrane, including the microvilli, in human
387 airway epithelium. Interestingly, about 25 years ago, prominin was shown to be localised to the apical
388 surface of various epithelial cells, where it is selectively associated with microvilli and microvillus-related
389 structures (Weigmann *et al.*, 1997). Furthermore, Weigmann and colleagues showed that prominin
390 expressed ectopically in non-epithelial cells was also selectively found in microvillus-like protrusions of the
391 plasma membrane. Two years later, (Corbeil *et al.*, 1999) showed that prominin contains dual targeting
392 information, for direct delivery to the apical domain of the plasma membrane and for enrichment in the
393 microvilli subdomain. Furthermore, they showed that this dual targeting does not require the cytoplasmic C-
394 terminal tail of prominin (*i.e.*, part 814-865 of CD133). From the above results, it is tempting to assume that
395 the prominin-like part of Spike S is involved in the delivery of the virus to the apical domain of the plasma
396 membrane where the ACE2 proteins are located. This hypothesis is all the more tempting as the similarity
397 between Spike S and prominin does not concern the C-terminal part of prominin, which, as we have pointed
398 out above, is not necessary for prominin targeting (recall that we have shown that the 186-482 part of
399 human prominin-1 is similar to the 908-1254 part of Spike S). Unfortunately, to date, the molecular nature of
400 the prominin apical sorting signal is unknown. It has been suggested in (Weigmann *et al.*, 1997) that
401 prominin may interact with the actin cytoskeleton, or that plasma membrane protrusions may have a specific
402 lipid composition/organisation for which prominins may have a preference.

403 Finally, it should be noted that the "SARS-CoV(-1)" glycoprotein Spike, which, like SARS-CoV-2 Spike,
404 binds to human ACE2 (Li *et al.*, 2003), is also similar to human prominin-1. Specifically, using HHpred, we
405 showed that the 177-473 part of the latter is similar to the 890-1236 part of Spike (with an associated
406 probability of 0.95 – see Supplemental file 4, raw HHpred data). In contrast, the MERS-CoV Spike
407 glycoprotein (like SARS-CoV and SARS-CoV-2, MERS-CoV is a betacoronavirus), which uses human DPP4 as an
408 entry receptor (Raj *et al.*, 2013), is similar to human mucin-1: the 292-421 part of mucin-1 is similar – with an
409 associated probability of 0.89 – to the 1230-1344 part of MERS-CoV Spike (see Supplemental file 4, raw
410 HHpred data). It is also interesting to note that (Kotani *et al.*, 2022) showed that the DPP4 protein also
411 colocalises with ACE2 and CD133 in Caco-2 cells. This suggests that it is likely that *i/* different coronaviruses
412 compete at the same positions on the cell, but *ii/* use different entry receptors and therefore different types
413 of spike proteins to reach these sites and fuse with the cells.

414 **A previously unknown "quite robust" similarity was detected on ORF3a protein**

415 The 41-183 part of ORF3a (275 A.A.) shows similarities to human G Protein-Coupled Receptors (GPCRs)
416 (which are cell surface receptor proteins that detect molecules from outside the cell and trigger cellular
417 responses (Lagerström & Schiöth, 2008)) and in particular to the GPCRs annotated with the Pfam domain "7
418 transmembrane receptor (rhodopsin family)/7tm_1" (see Results section and Supplemental file 2). According
419 to Pfam, this family contains, among other GPCRs, members of the opsin family, which are considered typical
420 members of the rhodopsin superfamily.

421 The ORF3a protein of "SARS-CoV(-1)" has been shown to form an ion channel (Lu *et al.*, 2006). Recently,
422 (Kern *et al.*, 2021) presented Cryo-EM determined structures of SARS-CoV-2 ORF3a at a resolution of 2.1Å.
423 The authors provide evidence suggesting that ORF3a forms a large polar cavity in the inner half of the
424 transmembrane region (TM) that could form ionic conduction paths (TM1 (43-61), TM2 (68-99) and TM3
425 (103-133)). Interestingly, the similarity we detected on ORF3a (41-183) encompasses the transmembrane
426 portion of ORF3a (43-133) which could form ionic permeation pathways. As mentioned earlier, we have
427 shown that this part of ORF3a resembles many GPCRs which belong to the Rhodopsin family (22 of 28 human
428 proteins sharing similarities with ORF3a, see Supplemental file 2 for more details). It is interesting to note
429 that some GPCRs, called "Rhodopsin channels", directly form ion channels (see (Nagel *et al.*, 2002) and
430 (Nagel *et al.*, 2003)). From this point of view, our prediction is therefore in line with the work of (Kern *et al.*,
431 2021). However, it is worth mentioning that a recent work challenges the results of both (Kern *et al.*, 2021)
432 and (Lu *et al.*, 2006): (Miller *et al.*, 2023) provide evidence suggesting that while a narrow cavity is detected
433 in the SARS-CoV-2 ORF3a transmembrane region, it likely does not represent a functional ion-conducting
434 pore (the same holds true for SARS-CoV-1 ORF3a).

435 However, Finally, it should be noted that if our method is applied to the ORF3a of SARS-CoV(-1), no
436 similarities are identified. More precisely, none of the similarities found by HHpred are significant, *i.e.* the
437 probability of the best hit is 0.72, which is below our threshold of 0.8; moreover, this best hit does not
438 correspond to a GPCR (see Supplemental file 4). This result may suggest a lack of sensitivity of HHpred. That
439 said, although HHpred is a fairly effective tool for detecting very distant homologies, not all similarities are
440 detectable. Furthermore, although the ORF3a of SARS-CoV(-1) and SARS-CoV-2 share 72% sequence identity
441 and are similar in the arrangement of the TM domains, the differences observed in the ion channel
442 properties between these two proteins suggest a different mode of action between them (Zhang *et al.*,
443 2022).

444 Autoantibodies targeting GPCRs have been found in patients with COVID-19 and Long COVID-19. It is
445 therefore tempting to speculate that the similarity between ORF3a and certain human GPCRs could be the
446 cause of the autoimmune reactions observed. There is some evidence to support this hypothesis:

447 –Autoantibodies targeting GPCRs (and RAS-related molecules) have been shown to be associated with
448 the severity of COVID-19 (Cabral Marques *et al.*, 2022). Among the anti-GPCR autoantibodies, the authors of

449 the latter paper identified the chemokine receptor CXCR3 and the RAS-related molecule AGTR1 as antibody
450 targets with the strongest association with disease severity. Strikingly, of the 26 GPCRs we identified as
451 sharing similarity with ORF3a (see Supplemental file 2), 5 are chemokine receptors, including the chemokine
452 receptor CXCR3;

453 —Functional autoantibodies against G-protein-coupled receptors have been found in patients with
454 persistent symptoms of Long COVID-19 (Wallukat *et al.*, 2021). In particular, the authors of the latter paper
455 identified functional autoantibodies against the M2 muscarinic receptor in the blood of Long COVID patients.
456 Strikingly, of the 26 GPCRs we identified as sharing similarity with ORF3a (see Supplemental file 2), 3 are
457 muscarinic receptors, including the muscarinic acetylcholine M2 receptor. In the same study, functional
458 autoantibodies against the alpha-1 adrenoceptor and the beta-2 adrenoceptor were also identified.
459 Interestingly, of the 26 above-mentioned GPCRs, 3 are adrenoceptors, namely alpha-1D, alpha-2A and
460 alpha-2C (see Supplemental file 2).

461 To conclude this section, we would like to emphasize that our main goal is to identify similarities between
462 SARS-CoV-2 proteins and human proteins in order to gain a better understanding of the functions of SARS-
463 CoV-2 proteins, rather than seeking mimics that could trigger autoimmune processes. This problem is usually
464 solved by searching for n-mers, which is obviously not done here. (Khavinson *et al.*, 2021), for example,
465 specifically addresses this problem and concludes that ORF3a does not appear to be involved in triggering an
466 autoimmune response. Furthermore, based solely on the similarity between ORF3a and certain human
467 GPCRs targeted by autoantibodies in patients with COVID-19 and Long COVID-19, it is difficult to state that
468 this similarity is the cause of the autoimmune phenomena observed. As Cabral Marques *et al.* (2022) point
469 out, the mechanisms by which SARS-CoV-2 infection triggers the production of autoantibodies remain
470 unknown to this day; according to these authors, molecular mimicry between SARS-CoV-2 and certain human
471 proteins is obviously not the only hypothesis to explain these phenomena: a hyperinflammatory response
472 triggered by the virus could cause tissue damage, leading to systemic autoimmune reactions. However, our
473 results suggest that further studies should be conducted.

474 **Comparison of our results with those of "Pfam clans"**

475 As indicated in the introduction to this article (see also Supplemental file 1), of the 40 Pfam domains that
476 annotate SARS-CoV-2 proteins, only one domain is not confined to viruses, the Macro domain that annotates
477 NSP3. This observation can be modulated at the level of Pfam clans which are collections of related domains.
478 At this level, 12 domains belong to clans whose domains are not strictly viral (see Supplemental file 1). These
479 clans allow the annotation of the following 9 proteins (more generally, of only part of each protein): NSP3,
480 NSP5, NSP13, NSP14, NSP15, NSP16, ORF7a, ORF8, and Spike S. 4 of these proteins are annotated by both
481 Pfam and our approach: NSP3, NSP13, NSP16 and Spike S. In the case of NSP3, NSP13 and NSP16, the
482 annotations are similar (note however that for NSP3, Pfam detects two domains related to the MACRO clan;
483 only one Macro domain is detected by our approach) whereas in the case of Spike S, our annotations refer to
484 a different part of the protein than that annotated by Pfam. We also identified similarities, not restricted to
485 viruses unlike Pfam, for ORF3a and NSP2.

486 **Evaluation of our results in light of the known weaknesses of HHpred**

487 As reported in (Gabler *et al.*, 2020) and (Kuchibhatla *et al.*, 2014), some false positive HHpred hits may
488 have high scores because they have coiled-coil, transmembrane or low complexity segments. Of our 6
489 "robust similarities", 2 have transmembrane segments and/or disordered areas (according to InterPro
490 annotations).

491 ORF3a

492 As previously indicated, ORF3a shares similarity with G Protein-Coupled Receptors (GPCRs) annotated
493 with the Pfam domains "7 transmembrane receptor (rhodopsin family)/7tm_1" or "7 transmembrane
494 receptor (secretin-like) 7tm_2" (see Results or Supplemental file 2).

495 Since transmembrane proteins are a large family of proteins – according to UniProt, out of 80581
496 proteins expressed by humans, 13876 are transmembrane proteins – it is legitimate to ask whether the
497 (observed) distribution of transmembrane proteins found by HHpred – out of 28 proteins found by HHpred,
498 28 are transmembrane proteins – is the same as the (expected) distribution of transmembrane proteins in
499 UniProt. Using a Fisher's exact test, we conclude (see Supplemental file 3 for proof) that the results found by
500 HHpred are not randomly drawn from the UniProt human proteome (p-value = 6.2059249716913E-11).

501 Similarly, as transmembrane proteins can be grouped into many different classes (the Pfam clan "Family
502 A G protein-coupled receptor-like superfamily", to which 7tm_1 and 7tm_2 belong, alone contains 53
503 different domains), it can also be argued that the similarities found by HHpred are due to chance. Of the 28
504 transmembrane proteins found by HHpred, 26 belong to the 7tm_1 or 7tm_2 classes. Knowing that the
505 number of human proteins belonging to the 7tm_1 or 7tm_2 classes is – according to UniProt – 540, we
506 show (see Supplemental file 3 for proof) using a Fisher's exact test that the results obtained by HHpred do
507 not arise from random selection within the different classes of the transmembrane protein family (p-value =
508 2.8739559680731E-12).

509 Spike glycoprotein

510 As shown previously, the 908-1254 part of the Spike S protein of SARS-CoV-2 is similar to the 186-482
511 part of human prominin-1. The 179-432 part of this prominin is annotated as
512 "NON_CYTOPLASMIC_DOMAIN" (*i.e.* non-cytoplasmic loops of a TM protein) by Phobius (for completeness,
513 note that the 253-283 part is annotated as a coil by COILS).

514 In contrast to the case of ORF3a, no reliable statistical test can be performed here (the number of human
515 prominins, *i.e.* proteins annotated by Pfam as "prominin" (Pfam PF05478), is 5). However, such a calculation
516 seems unnecessary here. HHpred identified a similarity between Spike S and human and fly prominins (see
517 Results section). Human and fly belonging to lineages that were separated over 700 million years ago
518 (median time of divergence 694 MYA (see <http://timetree.org/>, (Kumar et al., 2017)), this similarity is clearly
519 not a coincidence (unless one imagines a recent horizontal transfer).

520

Conclusion

521 We used HHpred to search for similarities between SARS-Cov-2 and human proteins. To avoid false
522 predictions, the robustness of each similarity was assessed using a procedure based on "test
523 sets/proteomes". We found six robust similarities in six different proteins, of which three are already
524 documented, one is in agreement with recent crystallographic results, and two are not reported in the
525 literature. We focused on these last two similarities and showed how they open new avenues of research to
526 better understand this virus. Obviously, our work is limited to making predictions that need to be validated
527 experimentally. Furthermore, the origin of the similarities (evolutionary convergence, horizontal transfer,
528 etc.) has not been addressed in this work. Nevertheless, we believe that our approach (or one similar to it)
529 can be profitably used to open up lines of research and to improve the annotation of any virus, especially
530 "orphan viruses", *i.e.* viruses which, for various reasons, are far much less studied than SARS-CoV-2.

531

Acknowledgements

532 We would like to thank all those who initiated Pfam, InterPro and HH-suite and all those who maintain and
533 improve these databases and tools for the benefit of our community. Without them, the work presented
534 here would not have been possible. We also warmly thank Joël Pothier (ISYEB, MNHN) whose comments
535 have clearly contributed to improve the quality of this paper.

536

Supplementary information availability

537 Supplemental files 1, 2, 3 and 4 are available online (*cf.* bioRxiv, DOI of this document).

538

Conflict of interest disclosure

539 The author declares that he complies with the PCI rule of having no financial conflicts of interest in relation
540 to the content of the article.

541

References

- 542 Boulant S., Stanifer M., Lozach P-Y. Dynamics of virus-receptor interactions in virus binding, signaling, and
543 endocytosis. *Viruses*. 2015 Jun 2;7(6):2794-815. <https://doi.org/10.3390/v7062747>
- 544 Blum M., Chang H., Chuguransky S., Grego T., Kandasaamy S., Mitchell A., Nuka G., Paysan-Lafosse T., Qureshi
545 M., Raj S., Richardson L., Salazar G. A., Williams L., Bork P., Bridge A., Gough J., Haft D. H., Letunic I.,
546 Marchler-Bauer A., Mi H., Natale D. A., Necci M., Orengo C. A., Pandurangan A. P., Rivoire C., Sigrist C. J.
547 A., Sillitoe I., Thanki N., Thomas P. D., Tosatto S. C. E., Wu C. H., Bateman A. and Finn R. D. The InterPro
548 protein families and domains database: 20 years on. *Nucleic Acids Research*, Nov 2020.
549 <https://doi.org/10.1093/nar/gkaa977>
- 550
- 551 [Cabral-Marques O., Halpert G., Schimke L. F., Ostrinski Y., Vojdani A., Baiocchi G. C., Freire P. P., Filgueiras I.](#)
552 [S., Zyskind I., Lattin M. T., Tran F., Schreiber S., Marques A. H. C., Praça D. R., Fonseca D. L. M., Humrich J. Y.,](#)
553 [Müller A., Giiil L. M., Graßhoff H., Schumann A., Hackel A., Junker J., Meyer C., Ochs H. D., Lavi Y. B.,](#)
554 [Scheibebogen C., Dechend R., Jurisica I., Schulze-Forster K., Silverberg J. I., Amital H., Zimmerman J.,](#)
555 [Heidecke H., Rosenberg A. Z., Riemekasten G., and Shoenfeld Y. Autoantibodies targeting GPCRs and RAS-](#)
556 [related molecules associate with COVID-19 severity. *Nature Communications*. 2022 Mar 9;13\(1\):1220.](#)
557 <https://doi.org/10.1038/s41467-022-28905-5>
- 558 Corbeil D., Röper K., Hannah M. J., Hellwig A., Huttner W. B. Selective localization of the polytopic membrane
559 protein prominin in microvilli of epithelial cells - a combination of apical sorting and retention in plasma
560 membrane protrusions. *J Cell Sci* (1999) 112 (7): 1023-1033. <https://doi.org/10.1242/jcs.112.7.1023>
- 561 Elde N. C., Malik H. S. The evolutionary conundrum of pathogen mimicry. *Nature Reviews Microbiology*. 2009
562 Nov;7(11):787-97. <https://doi.org/10.1038/nrmicro2222>
- 563 Forni D., Cagliani R., Molteni C., Arrigoni F., Mozzi A., Clerici M., De Gioia L., Sironi M. Homology-based
564 classification of accessory proteins in coronavirus genomes uncovers extremely dynamic evolution of
565 gene content. *Molecular Ecology*. 2022 Jul;31(13):3672-3692. <https://doi.org/10.1111/mec.16531>
- 566 Gabler F., Nam S.-Z., Till S., Mirdita M., Steinegger M., Söding J., Lupas A. N., Alva V. Protein Sequence
567 Analysis Using the MPI Bioinformatics Toolkit. *Current Protocols in Bioinformatics*. 2020 Dec;72(1).
568 <https://doi.org/10.1093/nar/gkl217>
- 569 Hoffmann M., Kleine-Weber H., Schroeder S., Krüger N., Herrler T., Erichsen S., Schiergens T. S., Herrler G.,
570 Wu N.-H., Nitsche A., Müller M. A., Drosten C., Pöhlmann S. SARS-CoV-2 Cell Entry Depends on ACE2 and

571 TMPRSS2 and Is Blocked by a Clinically Proven Protease Inhibitor. *Cell*. 2020 Apr 16;181(2):271-280.
572 <https://doi.org/10.1016/j.cell.2020.02.052>

573 Huang Y., Yang C., Xu X.-F., Xu W., Liu S.-W. Structural and functional properties of SARS-CoV-2 spike protein:
574 potential antiviral drug development for COVID-19. *Acta Pharmacologica Sinica*. 2020 Sep;41(9):1141-1149.
575 <https://doi.org/10.1038/s41401-020-0485-4>

576 Kern D. M., Sorum B., Mali S. S., Hoel C. M., Sridharan S., Remis J. P., Toso D. B., Kotecha A., Bautista D. M.,
577 Brohawn S. G. Cryo-EM structure of SARS-CoV-2 ORF3a in lipid nanodiscs. *Nat Struct Mol Bio*. 2021
578 Jul;28(7):573-582. <https://doi.org/10.1038/s41594-021-00619-0>

579 ~~Khavinson V., Terekhov A., Kormilets D., Maryanovich A. Homology between SARS-CoV-2 and human~~
580 ~~proteins. *Sci Rep*. 2021 Aug 25;11(1):17199. <https://doi.org/10.1038/s41598-021-96233-7>~~

581 Kotani N., Nakano T., Kuwahara R. Host cell membrane proteins located near SARS-CoV-2 spike protein
582 attachment sites are identified using proximity labeling and proteomic analysis. *J Biol Chem*. 2022
583 Nov;298(11):102500. <https://doi.org/10.1016/j.jbc.2022.102500>

584 Kuchibhatla D. B., Sherman W. A., Chung B. Y. W., Cook S., Schneider G., Eisenhaber B., Karlin D. G. Powerful
585 sequence similarity search methods and in-depth manual analyses can identify remote homologs in many
586 apparently "orphan" viral proteins. *Journal of Virology*, 2014 Jan; 88(1): 10-20.
587 <https://doi.org/10.1128/JVI.02595-13>

588 Kumar S., Stecher G., Suleski M., Hedges S. B. TimeTree: A Resource for Timelines, Timetrees, and Divergence
589 Times. *Mol Biol Evol*. 2017 Jul 1;34(7):1812-1819. <https://doi.org/10.1093/molbev/msx116>

590 Lagerström M. C., Schiöth H. B. Structural diversity of G protein-coupled receptors and significance for drug
591 discovery. *Nature Reviews Drug Discovery*. 2008 Apr;7(4):339-57. <https://doi.org/10.1038/nrd2518>

592 Li W., Moore M. J., Vasilieva N., Sui J., Wong S. K., Berne M. A., Somasundaran M., Sullivan J., Luzuriaga K.,
593 Greenough T. C., Choe H., Farzan M. Angiotensin-converting enzyme 2 is a functional receptor for the
594 SARS coronavirus. *Nature*. 2003 Nov 27;426(6965):450-4. <https://doi.org/10.1038/nature02145>

595 Lu W., Zheng B.J., Xu K., Schwarz W., Du L., Wong C.K., Chen J., Duan S., Deubel V., Sun B. Severe acute
596 respiratory syndrome-associated coronavirus 3a protein forms an ion channel and modulates virus
597 release. *Proc. Natl. Acad. Sci. USA*. 2006;103:12540-12545. <https://doi.org/10.1073/pnas.0605402103>

598 Ma J., Chen Y., Wu W., Chen Z. Structure and Function of N-Terminal Zinc Finger Domain of SARS-CoV-2
599 NSP2. *Virologica sinica*. 2021 Oct;36(5):1104-1112. <https://doi.org/10.1007/s12250-021-00431-6>

600 ~~Miller A. N., Houlihan P. R., Matamala E., Cabezas-Bratesco D., Young Lee G., Cristofori-Armstrong B., Dilan T.~~
601 ~~L., Sanchez-Martinez S., Matthies D., Yan R., Yu Z., Ren D., Brauchi S. E., Clapham D. E. The SARS-CoV-2~~
602 ~~accessory protein Orf3a is not an ion channel, but does interact with trafficking proteins. *Elife*. 2023 Jan~~
603 ~~25;12:e84477. <https://elifesciences.org/articles/84477>~~

604 Mistry J., Chuguransky S., Williams L., Qureshi M., Salazar G. A., Sonnhammer E. L. L., Tosatto S. C. E., Paladin
605 L., Raj S., Richardson L. J., Finn R. D., Bateman A. Pfam: The protein families database in 2021. *Nucleic*
606 *Acids Research*. 2021 8 January Volume 49, Issue D1, Pages D412-D419.
607 <https://doi.org/10.1093/nar/gkaa913>

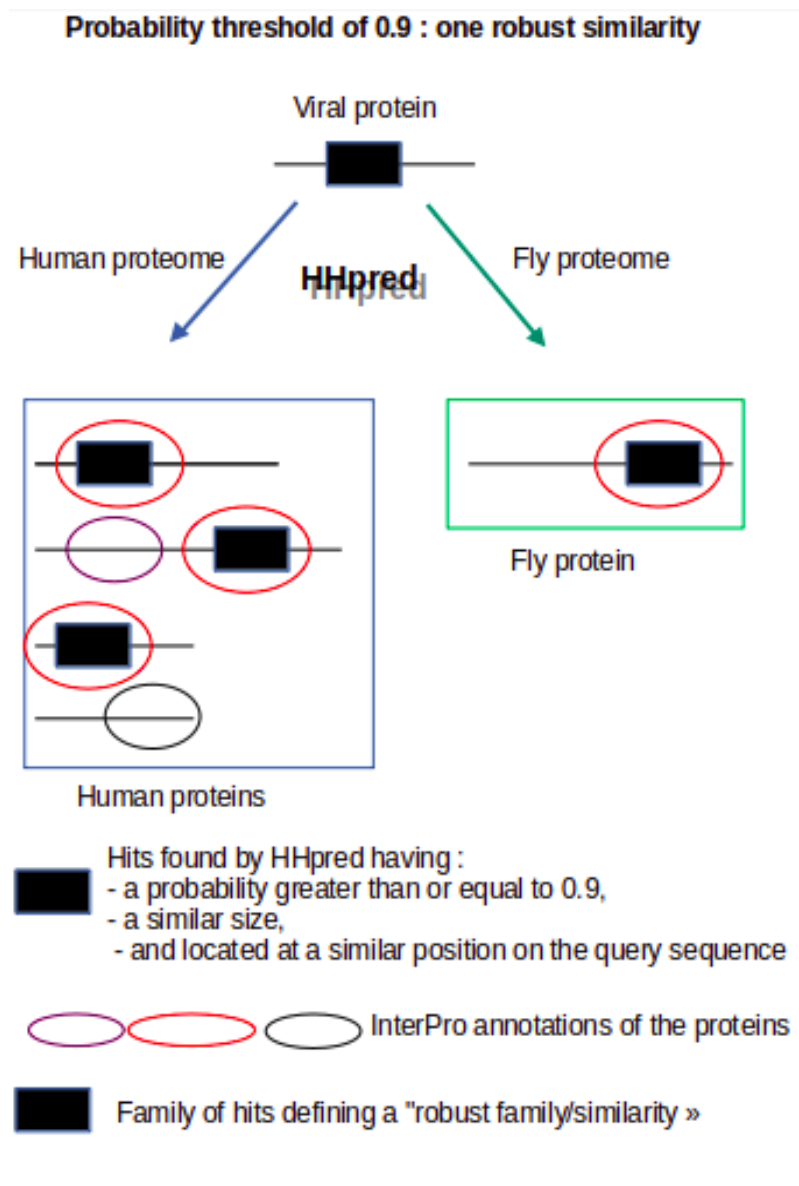
608 Nagel G., Ollig D., Fuhrmann M., Kateriya S., Musti A. M., Bamberg E., Hegemann P. Channelrhodopsin-1: a
609 light-gated proton channel in green algae. *Science*. 2002 Jun 28;296(5577):2395-8.
610 <https://doi.org/10.1126/science.1072068>

611 Nagel G., Szellas T., Huhn W., Kateriya S., Adeishvili N., Berthold P., Ollig D., Hegemann P., and Bamberg E.
612 Channelrhodopsin-2, a directly light-gated cation-selective membrane channel. *Proc Natl Acad Sci U S A*.
613 2003 Nov 25; 100(24): 13940-13945. <https://doi.org/10.1073/pnas.1936192100>

614 O'Donoghue S., Schafferhans A., Sikta N., Stolte C., Kaur S., Ho B. K., Anderson S., Procter J. B., Dallago C.,
615 Bordin N., Adcock M., Rost B. SARS-CoV-2 structural coverage map reveals viral protein assembly,
616 mimicry, and hijacking mechanisms. *Molecular System Biology*. 2021 Sep;17(9).
617 <https://doi.org/10.15252/msb.202010079>

618 Pinto A. L., Rai R. K., Brown J. C., Griffin P., Edgar J. R., Shah A., Singanayagam A., Hogg C., Barclay W. S.,
619 Futter C. E., Burgoyne T. Ultrastructural insight into SARS-CoV-2 entry and budding in human airway
620 epithelium. *Nat Comm.* 2022 Mar 25;13(1):1609. <https://doi.org/10.1038/s41467-022-29255-y>
621 Raj V. S., Mou H., Smits S. L., Dekkers D. H. W., Müller M. A., Dijkman R., Muth D., Demmers J. A. A., Zaki A.,
622 Fouchier Ron A. M., Thiel V., Drosten C., Rottier P. J. M., Osterhaus A. D. M. E., Bosch B. J., Haagmans B. L.
623 Dipeptidyl peptidase 4 is a functional receptor for the emerging human coronavirus-EMC. *Nature.*
624 2013;495:251–254. <https://doi.org/10.1038/nature12005>
625 Rouaud F., Méan I., Citi S. The ACE2 Receptor for Coronavirus Entry Is Localized at Apical Cell-Cell Junctions of
626 Epithelial Cells. *Cells.* 2022 Feb 11;11(4):627. <https://doi.org/10.3390/cells11040627>
627 Terrapon N. , Gascuel O., Maréchal E., Bréhélin L. Fitting hidden Markov models of protein domains to a
628 target species: application to *Plasmodium falciparum*. *BMC Bioinformatics.* 2012 May 1;13:67.
629 <https://doi.org/10.1186/1471-2105-13-67>
630 UniProt Consortium. UniProt: the universal protein knowledgebase in 2021. *Nucleic Acids Research.* 2021 Jan
631 8;49(D1):D480-D489. <https://doi.org/10.1093/nar/gkaa1100>
632 ~~Wallukat G., Hohberger B., Wenzel K., Fürst J., Schulze-Rothe S., Wallukat A., Hönicke A. S., Müller J.~~
633 ~~Functional autoantibodies against G-protein-coupled receptors in patients with persistent Long-COVID-19~~
634 ~~symptoms. Journal of Translational Autoimmunity. 2021;4:100100.~~
635 ~~<https://doi.org/10.1016/j.jtauto.2021.100100>~~
636 Weigmann A., Corbeil D., Hellwig A., Huttner W. B. Prominin, a novel microvilli-specific polytopic membrane
637 protein of the apical surface of epithelial cells, is targeted to plasmalemmal protrusions of non-epithelial
638 cells. *PNAS.* 1997 Nov 11;94(23):12425-30. <https://doi.org/10.1073/pnas.94.23.12425>
639 Xia S., Zhu Y., Liu M., Lan Q., Xu W., Wu Y., Ying T., Liu S., Shi Z., Jiang S., Lu L. Fusion mechanism of 2019-
640 nCoV and fusion inhibitors targeting HR1 domain in spike protein. *Cell Mol Immunol.* 2020 Jul;17(7):765-
641 767. <https://doi.org/10.1038/s41423-020-0374-2>
642 Zhang J., Ejikemeuwa A., Gerzanich V., Nasr M., Tang Q., Simard J.M., Zhao R.Y. Understanding the Role of
643 SARS-CoV-2 ORF3a in Viral Pathogenesis and COVID-19. *Front Microbiol.* 2022 Mar 9;13:854567.
644 <https://doi.org/10.3389/fmicb.2022.854567>.

645
 646
 647
 648
 649
 650
 651
 652
 653
 654
 655
 656
 657
 658
 659
 660
 661
 662
 663
 664
 665
 666
 667
 668
 669
 670
 671
 672
 673
 674
 675
 676
 677
 678
 679



680
 681
 682
 683
 684
 685
 686
 687
 688
 689

Figure - Using HHpred, a viral protein is compared to the human proteome and to a set of other proteomes, called "test" proteomes, which include the fly proteome. The probability threshold was set at 0.9, so only hits with a probability value of 0.90 or greater are considered relevant here. 4 homologous hits (*i.e.*, hits of similar sizes and located at a similar position on the query sequence) exceeding the given threshold were found by HHpred (black boxes): 3 are found in humans and one in flies; the InterPro annotation of all the "black box" hits are the same (red oval); as the annotations of all these homologous hits are identical and at least one of these hits belongs to a test proteome, the corresponding family of homologous hits is considered to be a "robust/similar family"; this similarity will be used to annotate the corresponding hit on the viral protein.