

Characterization of accessory genes in coronavirus genomes

Christian Jean Michel

Laboratoire ICube

Claudine Mayer

Institut Pasteur

Olivier Poch

Laboratoire ICube

Julie Dawn Thompson (✉ thompson@unistra.fr)

Laboratoire ICube <https://orcid.org/0000-0003-4893-3478>

Research

Keywords: COVID-19; SARS-CoV-2; SARS-CoV; coronavirus; accessory genes; ORF prediction; circular code motifs

Posted Date: June 9th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-32190/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published on August 27th, 2020. See the published version at <https://doi.org/10.1186/s12985-020-01402-1>.

1 **Characterization of accessory genes in coronavirus genomes**

2 Christian Jean Michel¹, Claudine Mayer^{1,2,3}, Olivier Poch¹ and Julie Dawn Thompson^{1,*}

3
4 *¹ Department of Computer Science, ICube, CNRS, University of Strasbourg, Strasbourg, F-67412,*
5 *France*

6 *² Unité de Microbiologie Structurale, Institut Pasteur, CNRS UMR 3528, 75724 Paris Cedex 15, France*

7 *³ Université Paris Diderot, Sorbonne Paris Cité, 75724 Paris Cedex 15, France*

8 Christian Jean Michel: c.michel@unistra.fr

9 Claudine Mayer: mayer@pasteur.fr

10 Olivier Poch: olivier.poch@unistra.fr

11 Julie Dawn Thompson: thompson@unistra.fr

12

13 ***Corresponding author:**

14 Name: Julie D. Thompson

15 Address: Department of Computer Science, ICube, Strasbourg, F-67412, France

16 Phone: (33) 0368853296

17 Email: thompson@unistra.fr

18

19

20 **Abstract**

21 **Background:** The Covid19 infection is caused by the SARS-CoV-2 virus, a novel member of
22 the coronavirus (CoV) family. CoV genomes code for a ORF1a / ORF1ab polyprotein and four
23 structural proteins widely studied as major drug targets. The genomes also contain a variable
24 number of open reading frames (ORFs) coding for accessory proteins that are not essential for
25 virus replication, but appear to have a role in pathogenesis. The accessory proteins have been
26 less well characterized and are difficult to predict by classical bioinformatics methods.

27 **Methods:** We propose a computational tool GOFIX to characterize potential ORFs in virus
28 genomes. In particular, ORF coding potential is estimated by searching for enrichment in motifs
29 of the *X* circular code, that is known to be over-represented in the reading frames of viral genes.

30 **Results:** We applied GOFIX to study the SARS-CoV-2 and related genomes including SARS-
31 CoV and SARS-like viruses from bat, civet and pangolin hosts, focusing on the accessory
32 proteins. Our analysis provides evidence supporting the presence of overlapping ORFs 7b, 9b
33 and 9c in all the genomes and thus helps to resolve some differences in current genome
34 annotations. In contrast, we predict that ORF3b is not functional in all genomes. Novel putative
35 ORFs were also predicted, including a truncated form of the ORF10 previously identified in
36 SARS-CoV-2 and a little known ORF overlapping the Spike protein in Civet-CoV and SARS-
37 CoV.

38 **Conclusions:** Our findings contribute to characterizing sequence properties of accessory genes
39 of SARS coronaviruses, and especially the newly acquired genes making use of overlapping
40 reading frames.

41

42 **Keywords:** COVID-19; SARS-CoV-2; SARS-CoV; coronavirus; accessory genes; ORF
43 prediction; circular code motifs

44

46 **Background**

47 Coronaviruses (CoVs) cause respiratory and intestinal infections in animals and humans [1].
48 They were not considered to be highly pathogenic to humans until the last two decades, which
49 have seen three outbreaks of highly transmissible and pathogenic coronaviruses, including
50 SARS-CoV (severe acute respiratory syndrome coronavirus), MERS-CoV (Middle East
51 respiratory syndrome coronavirus), and SARS-CoV-2 (which causes the disease COVID-19).
52 Other human coronaviruses (such as HCoV-NL63, HCoV-229E, HCoV-OC43 or HKU1)
53 generally induce only mild upper respiratory diseases in immunocompetent hosts, although
54 some may cause severe infections in infants, young children and elderly individuals [1].

55 Extensive studies of human coronaviruses have led to a better understanding of coronavirus
56 biology. Coronaviruses belong to the family *Coronaviridae* in the order *nidovirales*. Whereas
57 MERS-CoV is a member of the *Merbecovirus* subgenus, phylogenetic analyses indicated that
58 SARS-CoV-2 clusters with SARS-CoV in the *Sarbecovirus* subgenus [2]. All human
59 coronaviruses are considered to have animal origins. SARS-CoV, MERS-CoV and SARS-
60 CoV-2 are assumed to have originated in bats [1]. It is widely believed that SARS-CoV and
61 SARS-CoV-2 were transmitted directly to humans from market civets and pangolin,
62 respectively, based on the sequence analyses of CoV isolated from these animals and from
63 infected patients.

64 All members of the coronavirus family are enveloped viruses that possess long positive-
65 sense, single-stranded RNA genomes ranging in size from 27–33 kb. The coronavirus genomes
66 encode five major open reading frames (ORFs), including a 5' frameshifted polyprotein
67 (ORF1a/ORF1ab) and four canonical 3' structural proteins, namely the spike (S), envelope (E),
68 membrane (M) and nucleocapsid (N) proteins, which are common to all coronaviruses [3]. In
69 addition, a number of subgroup-specific accessory genes are found interspersed among, or even
70 overlapping, the structural genes. Overlapping genes originate by a mechanism of overprinting,

71 in which nucleotide substitutions in a pre-existing frame induce the expression of a novel
72 protein in an alternative frame. The accessory proteins in coronaviruses vary in number,
73 location and size in the different viral subgroups, and are thought to contain additional functions
74 that are often not required for virus replication, but are involved in pathogenicity in the natural
75 host [4-5].

76 In the face of the ongoing COVID-19 pandemic, extensive worldwide research efforts have
77 focused on identifying coronavirus genetic variation and selection [6-8], in order to understand
78 the emergence of host/tissue specificities and to help develop efficient prevention and treatment
79 strategies. These studies are complemented by structural genomics [9-11], as well as
80 transcriptomics [12] and interactomics studies [13] of the structural and putative accessory
81 proteins.

82 However, there have been less studies of accessory proteins, for two main reasons [14]. First,
83 accessory proteins are often not essential for viral replication or structure, but play a role in
84 viral pathogenicity or spread by modulating the host interferon signaling pathways for example.
85 This has led to some contradictory experimental results concerning the presence or functionality
86 of accessory proteins. For instance, in a recent experiment [13] to characterize SARS-CoV-2
87 gene functions, 9 predicted accessory protein ORFs (3a, 3b, 6, 7a, 7b, 8, 9b, 9c, 10) were codon
88 optimized and successfully expressed in human cells, with the exception of ORF3b. However,
89 another recent study using DNA nanoball sequencing [12] concluded that the SARS-CoV-2
90 expresses only five canonical accessory ORFs (3a, 6, 7a, 7b, 8).

91 Second, bioinformatics approaches for the prediction of accessory proteins are challenged
92 by their complex nature as short, overlapping ORFs. Such proteins are known to have biased
93 amino acid sequences compared to non-overlapping proteins [15]. In addition, the homology-
94 based approaches widely used to predict ORFs in genomes are less useful here, because many
95 accessory proteins are lineage- or subgroup-specific. Thus, many state of the art viral genome

96 annotation systems, such as Vgas [16], only predict overlapping proteins if homology
97 information is available. Other methods have been developed dedicated specifically to the *ab*
98 *initio* prediction of overlapping genes, for example based on multiple sequence alignments and
99 statistical estimates of the degree of variability at synonymous sites [17] or sequence
100 simulations and calculation of expected ORF lengths [18].

101 Here, we propose a computational tool GOFIX (Gene prediction by Open reading Frame
102 Identification using *X* motifs) to predict potential ORFs in virus genomes. Using a complete
103 viral genome as input, GOFIX first locates all potential ORFs, defined as a region delineated
104 by start and stop codons. In order to predict functional ORFs, GOFIX calculates the enrichment
105 of the ORFs in *X* motifs, i.e. motifs of the *X* circular code [19], a set of 20 codons that are over-
106 represented in the reading frames of genes from a wide range of organisms. For example, in a
107 study of 299,401 genes from 5217 viruses [20] including double stranded and single stranded
108 DNA and RNA viruses, codons of the *X* circular code were found to occur preferentially in the
109 reading frame of the genes. This is an important property of viral genes, since it has been
110 suggested that *X* motifs at different locations in a gene may assist the ribosome to maintain and
111 synchronize the reading frame [21]. An initial evaluation test of the GOFIX method on a large
112 set of 80 virus genomes [15] showed that it achieves high sensitivity and specificity for the
113 prediction of experimentally verified overlapping proteins (manuscript in preparation). A major
114 advantage of our approach is that it requires only the sequence of the studied genome and does
115 not rely on any homology information. This allows us to detect novel ORFs that are specific to
116 a given lineage.

117 We applied GOFIX to study the SARS-CoV-2 genome and related SARS genomes, with a
118 main focus on the accessory proteins. Using the extensive experimental data concerning the
119 SARS-CoV genome and the expressed ORFs, we first show that the reading frames of the
120 SARS-CoV ORFs are enriched in *X* motifs, including most of the overlapping accessory

121 proteins. Exceptions include ORF3b and ORF8b which may not be functional. Then, we use
122 GOFIX to predict and compare putative genes in related genomes of SARS-like viruses from
123 bat, civet and pangolin hosts as well as human SARS-CoV-2.

124

125 **Methods**

126 *Genome sequences*

127 Viral genome sequences were downloaded from the Genbank database, as shown in Table
128 1. The Genbank reference genomes were used as representative genomes for SARS-CoV and
129 SARS-CoV-2. For the Bat-CoV, Civet-CoV and Pangolin-CoV genomes, we selected well
130 annotated Genbank entries having the highest number of annotated ORFs. All CDS annotations
131 were extracted from the Genbank files, and ORF names were standardized according to the
132 SARS-CoV-2 nomenclature (Table 2).

133

134 *Definition of X motif enrichment (XME) scores*

135 The X circular code contains the following 20 codons

$$X = \{AAC, AAT, ACC, ATC, ATT, CAG, CTC, CTG, GAA, GAC, \\ GAG, GAT, GCC, GGC, GGT, GTA, GTC, GTT, TAC, TTC\} \quad (1)$$

136 and has several strong mathematical properties [19]. In particular, it is self-complementary, i.e.
137 10 trinucleotides of X are complementary to the other 10 trinucleotides of X , and it is a circular
138 code. A circular code is defined as a set of words such that any motif obtained from this set,
139 allows to retrieve, maintain and synchronize the reading frame.

140 An X motif m is defined as a word containing only codons from the X circular code (1) with
141 length $|m| \geq 3$ codons and cardinality (i.e. number of unique codons) $c \geq 2$ codons. The
142 minimal length $|m| = 3$ codons was chosen based on a previous study showing that the
143 probability of retrieving the reading frame with an X motif of at least 3 codons is 99.9% [22].

144 The class of X motifs with cardinality $c < 2$ are excluded here because they are mostly
145 associated with the “pure” trinucleotide repeats often found in non-coding regions of genomes
146 [23].

147 The total length XL_f of all X motifs m_f of nucleotide length $|m_f|$ in a frame f (the reading
148 frame or one of the 2 shifted frames) of a nucleotide sequence s is defined as:

$$149 \quad XL_f = \sum_{m_f \in s} |m_f|.$$

150 Then the X motif enrichment XME_f in a frame f of a sequence s of nucleotide length l is
151 defined as:

$$152 \quad XME_f = \frac{100}{l_f} XL_f$$

153 where for non-overlapping ORFs: $l_f = l$, and for overlapping ORFs: $l_f = l - XL_g$ where
154 XL_g is the total length of all X motifs in the overlapped frame g .

155 Finally, for an ORF of length l and associated with a reading frame f , the X motif enrichment
156 score XME is defined as:

$$157 \quad XME = XME_f$$

158

159 *GOFIX method*

160 The GOFIX method will be described in detail in a separate manuscript. Briefly, the method
161 consists of two main steps:

- 162 (i) Identification of all potential ORFs. Using the complete genome sequences as input, all
163 potential ORFs in the positive sense are located, defined as a sequence region starting
164 with a start codon (AUG) and ending with a stop codon (UAA, UAG, UGA). For a
165 given region, if alternative start codons are found, the longest ORF is selected. In this
166 study, we selected all ORFs having a minimum length of 120 nucleotides (40 amino
167 acids).

168 (ii) Calculation of X motif enrichment scores. For each potential ORF, all X motifs in the
169 nucleotide sequence are identified in the three positive sense frames f using the
170 computational method described in [24]. For each identified potential ORF, the X motif
171 enrichment (XME_f and XME) scores are calculated as defined above. Based on our
172 benchmark studies (data not shown) of experimentally validated ORFs in a large set of
173 80 genomes [15], we set the threshold for prediction of a functional ORF to be $XME \geq$
174 5.

175

176 **Results**

177 *Initial study of SARS-CoV reference genome*

178 We first analyzed the complete genome of the well-studied SARS-CoV and plotted the X
179 motif enrichment (XME_f) scores calculated in a sliding window of 150 nucleotides for each of
180 the three positive sense frames (Fig. 1). We then mapped the ORF1ab, the four structural
181 proteins (S, E, M, N), and the nine generally accepted accessory genes (3a, 3b, 6, 7a, 7b, 8a,
182 8b, 9b, 9c) to the X enrichment plot.

183 We observe a tendency for the reading frames of the SARS-CoV ORFs to be enriched in X
184 motifs. For example, ORF1ab is the longest ORF, encoding a polyprotein, which is translated
185 by a -1 programmed ribosomal frameshift at position 13398. Sequences upstream and
186 downstream of the frameshift are enriched in X motifs in the corresponding reading frame
187 (green and yellow plots respectively in Fig. 1A). Other ORFs enriched in X motifs in the reading
188 frame include the S protein (yellow plot in Fig. 1B) and the E and M proteins (blue and green
189 plots respectively in Fig. 1C). The S, E and M ORFs are conserved in all coronavirus genomes
190 and code for structural proteins that together create the viral envelope.

191 The case of overlapping ORFs is more complex. For example, the last structural protein
192 coded by the N ORF is overlapped by two accessory genes: ORF9b and ORF9c. The sequence

193 regions containing the overlapping ORFs are characterized by an enrichment in X motifs in the
194 2 frames (green and blue plots in Fig. 1C).

195

196 *Characterization of known accessory genes in SARS-CoV*

197 The SARS-CoV genome is known to contain four structural proteins and nine accessory
198 proteins, namely ORFs 3a, 3b, 6, 7a, 7b, 8a, 8b, 9b and 9c. To verify that our approach can
199 predict the accessory genes in coronavirus genomes, we used GOFIX to identify all potential
200 ORFs in the complete SARS-CoV genome and calculate their X enrichment. Fig. 2 shows the
201 X motif enrichment (XME_f) scores calculated by GOFIX for the identified ORFs in the 3'
202 terminal region of the SARS-CoV genome.

203 The overall performance of GOFIX is shown in Table 3. Initially, GOFIX found 25 potential
204 ORFs (delineated by start and stop codons) in the 3' region (21492-29751) of SARS-CoV.
205 Twelve of these 25 potential ORFs were predicted to be non-functional (see Methods),
206 including 10 unknown ORFs mostly overlapping the S protein. Two previously annotated ORFs
207 were also predicted to be non-functional, namely ORF3b ($XME=1.9$) and ORF8b ($XME=0.0$)
208 that are discussed in detail below.

209 GOFIX predicts that 13 of the 25 potential ORFs are functional (with $XME>5$). These
210 include 11 previously annotated ORFs, namely S, 3a, E, M, 6, 7a, 7b, 8a, N, 9b, 9c. Two novel
211 ORFs are also predicted by the GOFIX method: ORF10 ($XME=15.8$) is located downstream of
212 the N gene (29415-29496) and a new ORF we called ORFSa ($XME=7.6$) that overlaps the S
213 gene (22732-22928). These novel ORFs are discussed in more detail below.

214

215 *Comparative analyses of accessory proteins in coronavirus genomes*

216 Having evaluated the GOFIX method on the well-studied SARS-CoV genome, we then used
217 it to characterize and compare the accessory proteins in representative strains of five

218 coronavirus genera, including SARS-CoV, SARS-CoV-2 and three viruses from animal hosts
219 with SARS-CoV-like infections. Bat is considered to be the most likely host origin of SARS-
220 CoV and SARS-CoV-2. It is generally considered that transmission to humans occurred *via* an
221 intermediate host. For SARS-CoV, civets probably acted as the intermediate host, while
222 pangolin has been proposed as the intermediate host in SARS-CoV-2 animal-to-human
223 transmission [25]. For each of the five genomes, we used GOFIX to predict all potential ORFs
224 in the complete genomes and calculated the *X* motif enrichment (XME) scores for each ORF.
225 Fig. 3 gives an overview of the predicted ORFs in each genome, confirming for example that
226 the structural proteins S, E, M and N, as well as the accessory proteins ORF6, ORF7a and
227 ORF7b are conserved and have XME scores above the defined threshold XME=5. However,
228 important differences in XME scores are observed for the remaining accessory protein ORFs.

229

230 *ORF3b may not code for a functional protein in all CoVs*

231 ORF3a codes for the largest accessory protein that comprises 274-275 amino acids (Fig. 4).
232 In SARS-CoV, ORF3a is not required for virus replication, but contributes to pathogenesis by
233 mediating trafficking of Spike (S protein) [4]. It is efficiently expressed on the cell surface, and
234 was easily detected in a majority of SARS patients. The XME scores for ORF3a in all the
235 genomes range from 13.8-19.3, *i.e.* almost 3 times greater than the defined threshold for
236 functional ORFs.

237 The ORF3b coding sequence overlaps the +1 reading frame of ORF3a and sometimes
238 extends beyond the start codon of the E gene. In SARS-CoV, it is proposed to antagonize
239 interferon (IFN) function by modulating the activity of IFN regulatory factor 3 (IRF3) [26].
240 However, immunohistochemical analyses of tissue biopsies and/or autopsies of SARS-CoV-
241 infected patients have failed to demonstrate the presence of ORF3b *in vivo*, and the presence of
242 ORF3b in SARS-CoV-infected Vero E6 cells is the only evidence for the expression of this

243 protein [27]. Furthermore, when mice are infected with mutant SARS-CoV lacking ORF3b, the
244 deletion viruses grow to levels similar to those of wild-type virus, which demonstrates that
245 SARS-CoV is able to inhibit the host IFN response without the 3b gene [28].

246 Bat-Cov and Civet-CoV also present ORF3b overlapping the 3' region of ORF3a (start
247 codon at nt 422), although the sequence of Bat-CoV ORF3b is shorter having a stop codon
248 within the ORF3a sequence (nt 764). We observe a single *X* motif in the ORF3b reading frame
249 of length 9 nucleotides (563-571), resulting in low XME scores of 2.6, 1.9 and 1.9 respectively
250 for Bat-CoV, Civet-CoV and SARS-CoV ORF3b. This ORF is not predicted to be present in
251 Pangolin-CoV or SARS-CoV-2 due to the introduction of a new stop codon (indicated by ***
252 in Fig. 4) and the loss of the *X* motif in the +1 reading frame.

253 However, a completely different ORF is identified in the Pangolin-CoV and SARS-CoV-2
254 sequences, overlapping the 5' region of ORF3a (132-305). This ORF is not annotated in the
255 SARS-CoV-2 reference genome (MT072688), but is annotated as ORF3b in the genome of
256 another SARS-CoV-2 strain isolated from the first U.S. case of COVID-19 (MN985325). The
257 Pangolin-CoV ORF3b sequence contains one *X* motif in the reading frame of length 9
258 nucleotides (183-191), but the *X* motif is lost in the SARS-CoV-2 genome.

259

260 *ORF8: a rapidly evolving region of SARS-CoV genomes*

261 Previously shown to be a recombination hotspot, ORF8 is one of the most rapidly evolving
262 regions among SARS-CoV genomes [29]. Furthermore, the evolution of ORF8 is supposed to
263 play a significant role in adaptation to the human host following interspecies transmission and
264 virus replicative efficiency [30].

265 In SARS-CoV isolated from bats and civets (as well as early human isolates of the SARS-
266 CoV outbreak in 2003: data not shown), ORF8 encodes a single protein of length 122 amino
267 acids (Fig. 5). However, in SARS-CoV isolated from humans during the peak of the epidemic,

268 there is a 29-nt deletion in the middle of ORF8, resulting in the splitting of ORF8 into two
269 smaller ORFs, namely ORF8a and ORF8b [31]. ORF8a and ORF8b encode a 39 amino acid
270 and 84 amino acid polypeptide, respectively. The XME scores in these ORFs are in line with
271 the known experimental evidence concerning their functions. ORF8a has an XME score of 15.3
272 in SARS-CoV and anti-p8a antibodies were identified in some patients with SARS [32]. In
273 contrast, ORF8b has no *X* motifs in the reading frame, and protein 8b was not detected in SARS-
274 CoV-infected Vero E6 cells [31].

275 It is interesting to note that although Civet-CoV has a full-length ORF8, it has a low XME
276 score (XME=4.9) compared to Bat-CoV (XME=9.9). Thus, it is tempting to suggest that the
277 loss of *X* motifs in transmission of the virus from bats to civets is somehow linked to the loss
278 of ORF8 in the transmission from civets to humans. Both Pangolin-CoV and most SARS-CoV-
279 2 strains contain the full length ORF8, with XME scores of 23.1 and 12.4 respectively.
280 However, a 382-nt deletion has been reported recently covering almost the entire ORF8 of
281 SARS-CoV-2 obtained from eight hospitalized patients in Singapore, that has been
282 hypothesized to lead to an attenuated phenotype of SARS-CoV-2 [33].

283

284 *Characterization of ORFs overlapping the N gene*

285 The annotation of functional ORFs overlapping the N gene is variable in the different
286 genomes studied here. In SARS-CoV, only ORF9b has been observed to be translated, probably
287 *via* a ribosomal leaky scanning mechanism and may have a function during virus assembly
288 [30,34]. ORF9b limits host cell interferon responses by targeting the mitochondrial-associated
289 adaptor molecule (MAVS) signalosome. However, some SARS-CoV strains have an additional
290 ORF9c, annotated as a hypothetical protein (e.g. Genbank:AY274119). For Bat-CoV and
291 Pangolin-CoV, no overlapping genes are annotated in the corresponding Genbank entries. In
292 contrast, the Civet-CoV genome is predicted to contain both overlapping genes, ORF9b and

293 ORF9c. Similarly, the annotation of overlapping ORFs for SARS-CoV-2 is different depending
294 on the strain: the reference strain has no overlapping ORFs of the N gene, while the U.S. strain
295 has ORF9b and ORF9c (see Methods). ORF9c is described as a short polypeptide (70 amino
296 acids) dispensable for viral replication, but there is no data yet providing evidence that the
297 protein is expressed during SARS-CoV-2 infection.

298 Here, we predict that ORF9b and ORF9c are present in all genomes as overlapping ORFs
299 within the N gene (Fig. 6). Furthermore, Pangolin-CoV may also have an additional ORF, that
300 we called ORF9d (XME=12.7), in the 3' region of the N gene.

301

302 *Origin and evolution of ORF10*

303 ORF10 is proposed as unique to SARS-CoV-2 [35] and codes for a peptide only 38 amino
304 acids long. There is no data yet providing evidence that the protein is expressed during SARS-
305 CoV-2 infection. Therefore, we wanted to investigate the potential origin of this protein. New
306 proteins in viruses can originate from existing proteins acquired through horizontal gene
307 transfer or through gene duplication for example, or can be generated *de novo*. To determine
308 whether homologs of ORF10 are present in the other coronavirus genomes, we relaxed the
309 GOFIX parameters used to predict functional ORFs, and set the minimum ORF length to 60
310 nucleotides. The predicted ORFs in the different genomes are shown in Fig. 7. The Pangolin-
311 CoV genome contains a full-length ORF10 with XME=10.4, compared to the SARS-CoV-2
312 ORF10 with XME=20.2. A truncated version of ORF10 coding 26 amino acids is also detected
313 in the Bat-CoV, Civet-CoV and SARS-CoV genomes, although this short ORF is probably not
314 functional. We suggest that the ORF10 of SARS-CoV-2 thus evolved *via* the mutation of a stop
315 codon (TAA) at nt 76 and the addition of a new *X* motif of length 15 nucleotides in the 3'
316 region.

317

318 *Novel ORF overlapping the S gene*

319 The GOFIX method predicts a novel ORF, that we called ORFSa, overlapping the RBD
320 (Receptor Binding Domain) of the S (Spike) ORF in SARS-CoV (XME=7.6) and Civet-CoV
321 (XME=7.6). ORFSa is found in the +1 frame and codes for a protein with 64 amino acids, as
322 shown in Fig. 8. As the ORFSa sequence was not present in the Bat-CoV reference genome,
323 we also searched for the ORF in the genomes of other Bat-CoV strains, and found one
324 occurrence (XME=6.5) in the strain WIV16 (Genbank:KT444582) (Fig. 8), another bat
325 coronavirus that is closely related to SARS-CoV [36].

326 To investigate whether the novel ORFSa might be a functional protein in SARS-CoV, we
327 used BlastP to search the Genbank database for matches to viral proteins. A significant hit was
328 obtained with a sequence identity of 100% to the protein AAR84376, described as “putative
329 transmembrane protein 2d” from the genome of SARS coronavirus strain ZJ01 (AY28632). To
330 further characterize this putative protein, the Phobius web site (phobius.sbc.su.se) was used to
331 predict transmembrane (TM) helices. Two potential TM helices of nearly twenty amino acids
332 (residues 6-28 and 42-62) were predicted with a small inter-TM endodomain. Thus, this
333 potential double-membrane spanning small protein might complement the set of already known
334 SARS-CoV membrane proteins, namely the Spike (S), membrane (M) and envelope (E)
335 proteins.

336

337 **Discussion**

338 Coronaviruses are complex genomes with high plasticity in terms of gene content. This
339 feature is thought to contribute to their ability to adapt to specific hosts and to facilitate host
340 shifts [1]. It is therefore essential to characterize the coding potential of coronavirus genomes.
341 Here, we used an *ab initio* approach to identify potential functional ORFs in the genomes of a
342 set of representative SARS or SARS-like coronaviruses. Our method allows comprehensive

343 annotation of all ORFs. Surprisingly, the calculation of *X* motif enrichment is also accurate for
344 the detection of overlapping genes, even though the codon usage and amino acid composition
345 of overlapping genes is known to be significantly different from non-overlapping genes [15].

346 We showed that the predictions made by the GOFIX method have high sensitivity and
347 specificity compared to the known functional ORFs in the well characterized SARS-CoV. For
348 example, the annotated ORFs that have been described previously as non-functional or
349 redundant, notably ORF3b and ORF8b, are not predicted to be functional by GOFIX. In
350 contrast, we identified a putative small ORF overlapping the RBD of the Spike protein in
351 SARS-CoV, that is conserved in Civet-CoV and Bat-CoV strain WIV16. Protein sequence
352 analysis predicts that this novel ORF codes for a double-membrane spanning protein.

353 We then used the method GOFIX to compare all putative ORFs in representative genomes,
354 and showed that most are conserved in all genomes, including the structural proteins (S, E, M
355 and N) and accessory proteins 3a, 6, 7a, 7b, 9b and 9c. However, a number of ORFs were
356 predicted to be non-functional, notably ORF8b in SARS-CoV and ORF3b in all genomes. We
357 also identified potential new ORFs, including ORF9d in Pangolin-CoV and ORF10 in all
358 genomes.

359 Concerning SARS-CoV-2, to date, the coding potential of SARS-CoV-2 remains partially
360 unknown, and distinct studies have provided different genome annotations [37-38]. Overall, the
361 genome of SARS-CoV-2 has 89% nucleotide identity with bat SARS-like-CoV (ZXC21) and
362 82% with that of human SARS-CoV [38]. Our analysis shows that the genome organization is
363 conserved, and in particular ORF9b and ORF9c are predicted to be expressed in SARS-CoV-2
364 genome. As expected, the structural proteins, S, E, M and N are conserved and have similar
365 XME scores. Here, we have shown that ORF3a, ORF6 and ORF9b in SARS-CoV-2 also have
366 similar XME scores to SARS-CoV.

367 Previously identified differences include some interferon antagonists and inflammasome
368 activators encoded by SARS-CoV that are not conserved in SARS-CoV-2, in particular ORF8
369 in SARS-CoV-2 and ORF8a,b in SARS-CoV, as well as the completely different ORF3b [14].
370 ORF3b has 0 *X* motifs in SARS-CoV-2 and expression was not observed in recent experiments
371 aimed at characterizing the functions of SARS-CoV-2 proteins [13]. ORF10 is supposed to be
372 unique to SARS-CoV-2, however it is also present in the Pangolin-CoV genome and its origin
373 can be traced back to the Bat-CoV, where a truncated ORF of 26 amino acids, also present in
374 the civet and human SARS-CoV genomes, can be found. Here, we observe that ORF7a, ORF7b
375 and ORF9c have reduced XME scores in SARS-CoV-2. It remains to be seen whether these
376 differences reflect functional divergences between SARS-CoV and SARS-CoV-2.

377 **Conclusions**

378 In summary, we have developed a computational method GOFIX to characterize potential
379 ORFs in virus genomes and applied the method to study the SARS-CoV-2 and related genomes.
380 Our analysis of ORF coding potential helps to resolve some differences in current genome
381 annotations. In addition, we suggest that some annotated ORFs may not be functional and
382 predict novel putative ORFs in some genomes. Our findings contribute to characterizing
383 sequence properties of accessory genes of SARS coronaviruses, and especially the newly
384 acquired genes making use of overlapping reading frames.

385

386 **Abbreviations**

387 CoV: Coronavirus

388 GOFIX: Gene prediction by Open reading Frame Identification using *X* motifs

389 MERS: Middle East Respiratory Syndrome

390 nt: nucleotide

391 ORF: Open Reading Frame

392 SARS: Severe Acute Respiratory Syndrome

393 XME: X Motif Enrichment

394 **Declarations**

395 *Ethics approval and consent to participate*

396 Not applicable

397 *Consent for publication*

398 Not applicable

399 *Availability of data and materials*

400 The datasets analysed during the current study are available in the Genbank viral genomes

401 database, <https://www.ncbi.nlm.nih.gov/genome/viruses/>

402 *Competing interests*

403 The authors declare that they have no competing interests

404 *Funding*

405 This work was supported by the Agence Nationale de la Recherche (BIPBIP: ANR-10-BINF-

406 03-02; ReNaBi-IFB: ANR-11-INBS-0013, ELIXIR-EXCELERATE: GA-676559) and

407 institute funds from the French Centre National de la Recherche Scientifique, and the

408 University of Strasbourg.

409 *Authors' contributions*

410 CJM, OP and JDT participated in the conceptualization of the work. CJM and JDT developed

411 the methods and analyzed the data. CM and OP analyzed and interpreted the data. All authors

412 contributed to writing the manuscript. All authors read and approved the final manuscript.

413 *Acknowledgements*

414 The authors would like to thank the BiGEst Bioinformatics Platform for assistance.

415

416 **References**

417

- 418 1. Cui J, Li F, Shi Z. Origin and evolution of pathogenic coronaviruses. *Nat Rev Microbiol*
419 2019;17:181-92.
- 420 2. Coronaviridae Study Group of the International Committee on Taxonomy of Viruses. The
421 species Severe acute respiratory syndrome-related coronavirus: classifying 2019-nCoV
422 and naming it SARS-CoV-2. *Nature Microbiology*. 2020;5:536-44.
- 423 3. Ashour HM, Elkhatib WF, Rahman MM, Elshabrawy HA. Insights into the Recent 2019
424 Novel Coronavirus (SARS-CoV-2) in Light of Past Human Coronavirus Outbreaks.
425 *Pathogens*. 2020;9 pii:E186.
- 426 4. Schaecher SR and Pekosz A. SARS Coronavirus Accessory Gene Expression and
427 Function. *Molecular Biology of the SARS-Coronavirus*. 2009;22:153-66.
- 428 5. Liu DX, Fung TS, Chong KK, Shukla A, Hilgenfeld R. Accessory proteins of SARS-CoV
429 and other coronaviruses. *Antiviral Res*. 2014;109:97-109.
- 430 6. Cagliani R, Forni D, Clerici M, Sironi M. Computational inference of selection underlying
431 the evolution of the novel coronavirus, SARS-CoV-2. *J Virol*. 2020;pii:JVI.00411-20.
- 432 7. Khailany RA, Safdar M, Ozaslan M. Genomic characterization of a novel SARS-CoV-2.
433 *Gene Rep*. 2020;Apr 16:100682.
- 434 8. Wang C, Liu Z, Chen Z, Huang X, Xu M, He T, Zhang Z. The establishment of reference
435 sequence for SARS-CoV-2 and variation analysis. *J Med Virol*. 2020;Mar 13.
- 436 9. Jin Z, Du X, Xu Y, et al. Structure of M^{pro} from COVID-19 virus and discovery of its
437 inhibitors. *Nature* 2020;Apr 9. doi: 10.1038/s41586-020-2223-y.
- 438 10. Hussain M, Jabeen N, Raza F, Shabbir S, Baig AA, Amanullah A, Aziz B. Structural
439 variations in human ACE2 may influence its binding with SARS-CoV-2 spike protein. *J*
440 *Med Virol*. 2020;Apr 6.

- 441 11. Srinivasan S, Cui H, Gao Z, Liu M, Lu S, Mkandawire W, Narykov O, Sun M, Korkin D.
442 Structural Genomics of SARS-CoV-2 Indicates Evolutionary Conserved Functional
443 Regions of Viral Proteins. *Viruses*. 2020;12.
- 444 12. Kim D, Lee JY, Yang JS, Kim JW, Kim VN, Chang H. The Architecture of SARS-CoV-2
445 Transcriptome. *Cell*. 2020;pii: S0092-8674(20)30406-2.
- 446 13. Gordon DE et al., A SARS-CoV-2-human protein–protein interaction map reveals drug
447 targets and potential drug repurposing. *Nature*. 2020 Apr 30. doi: 10.1038/s41586-020-
448 2286-9.
- 449 14. Yuen KS, Ye ZW, Fung SY, Chan CP, Jin DY. SARS-CoV-2 and COVID-19: The most
450 important research questions. *Cell Biosci*. 2020;10:40.
- 451 15. Pavesi A, Vianelli A, Chirico N, Bao Y, Blinkova O, Belshaw R, Firth A, Karlin D.
452 Overlapping genes and the proteins they encode differ significantly in their sequence
453 composition from non-overlapping genes. *PLoS One*. 2018;13:e0202513.
- 454 16. Zhang KY, Gao YZ, Du MZ, Liu S, Dong C, Guo FB. Vgas: A Viral Genome Annotation
455 System. *Front Microbiol*. 2019;10:184.
- 456 17. Firth AE. Mapping overlapping functional elements embedded within the protein-coding
457 regions of RNA viruses. *Nucleic Acids Res*. 2014;4220:12425-39.
- 458 18. Schlub TE, Buchmann JP, Holmes EC. A Simple Method to Detect Candidate Overlapping
459 Genes in Viruses Using Single Genome Sequences. *Mol Biol Evol*. 2018;35:2572-81.
- 460 19. Arquès DG, Michel CJ. A complementary circular code in the protein coding genes. *J*
461 *Theor Biol*. 1996;182:45-58.
- 462 20. Michel CJ. The Maximal C^3 Self-Complementary Trinucleotide Circular Code X in Genes
463 of Bacteria, Archaea, Eukaryotes, Plasmids and Viruses. *Life (Basel)*. 2017;7 pii: E20.
- 464 21. Dila G, Ripp R, Mayer C, Poch O, Michel CJ, Thompson JD. Circular code motifs in the
465 ribosome: a missing link in the evolution of translation? *RNA*. 2019;25:1714-30.

- 466 22. Michel CJ. Circular Code Motifs in Transfer and 16S Ribosomal RNAs: A Possible
467 Translation Code in Genes *Comput Biol Chem.* 2012;37:24-37.
- 468 23. El Soufi K, Michel CJ. Unitary circular code motifs in genomes of eukaryotes. *Biosystems*
469 2017;153:45-62.
- 470 24. El Soufi K, Michel CJ. Circular code motifs in genomes of eukaryotes. *J Theor Biol.*
471 2016;408:198-212.
- 472 25. Zhang T, Wu Q, Zhang Z. Probable Pangolin Origin of SARS-CoV-2 Associated with the
473 COVID-19 Outbreak. *Curr Biol.* 2020;30:1578.
- 474 26. Kopecky-Bromberg SA, Martínez-Sobrido L, Frieman M, Baric RA, Palese P. Severe
475 acute respiratory syndrome coronavirus open reading frame (ORF) 3b, ORF 6, and
476 nucleocapsid proteins function as interferon antagonists. *J Virol.* 2007;81:548-57.
- 477 27. McBride R, Fielding BC. The role of severe acute respiratory syndrome (SARS)-
478 coronavirus accessory proteins in virus pathogenesis. *Viruses.* 2012;4:2902-23.
- 479 28. Yount B, Roberts RS, Sims AC, Deming D, Frieman MB, Sparks J, Denison MR, Davis
480 N, Baric RS. Severe acute respiratory syndrome coronavirus group-specific open reading
481 frames encode nonessential functions for replication in cell cultures and mice. *J Virol.*
482 2005;79:14909-22.
- 483 29. Ceraolo C, Giorgi FM. Genomic variance of the 2019-nCoV coronavirus. *J Med Virol.*
484 2020;92:522-8.
- 485 30. Xu K, Zheng B-J, Zeng R, Lu W, Lin Y-P. Severe acute respiratory syndrome coronavirus
486 accessory protein 9b is a virion-associated protein. *Virology.* 2009;388:279-85.
- 487 31. Oostra M, de Haan CA, Rottier PJ. The 29-nucleotide deletion present in human but not in
488 animal severe acute respiratory syndrome coronaviruses disrupts the functional expression
489 of open reading frame 8. *J Virol.* 2007;81:13876-88.

- 490 32. Chen C-Y, Ping Y-H, Lee H-C, Chen K-H, Lee Y-M. Open reading frame 8a of the human
491 severe acute respiratory syndrome coronavirus not only promotes viral replication but also
492 induces apoptosis. *J. Infect. Dis.* 2007;196:405-15.
- 493 33. Su YCF, Anderson DE, Young BE, Zhu F, Linster M, Kalimuddin S, Low JGH, Yan Z,
494 Jayakumar J, Sun L, Yan GZ, Mendenhall IH, Leo Y-S, Lye DC, Wang L-F, Smith GJD.
495 Discovery of a 382-nt deletion during the early evolution of SARS-CoV-2. *bioRxiv*
496 2020.03.11.987222; doi: <https://doi.org/10.1101/2020.03.11.987222>.
- 497 34. Shukla A, Hilgenfeld R. Acquisition of new protein domains by coronaviruses: analysis of
498 overlapping genes coding for proteins N and 9b in SARS coronavirus. *Virus Genes.*
499 2015;50:29-38.
- 500 35. Wu F, Zhao S, Yu B, et al. A new coronavirus associated with human respiratory disease
501 in China. *Nature.* 2020;579:265-269.
- 502 36. Yang XL, Hu B, Wang B, Wang MN, Zhang Q, Zhang W, Wu LJ, Ge XY, Zhang YZ,
503 Daszak P, Wang LF, Shi ZL. Isolation and Characterization of a Novel Bat Coronavirus
504 Closely Related to the Direct Progenitor of Severe Acute Respiratory Syndrome
505 Coronavirus. *J Virol.* 2015;90:3253-6.
- 506 37. Zhou P, et al. A pneumonia outbreak associated with a new coronavirus of probable bat
507 origin. *Nature.* 2020;579:270-3.
- 508 38. Chan JF, Kok KH, Zhu Z, Chu H, To KK, Yuan S, Yuen KY. Genomic characterization of
509 the 2019 novel human-pathogenic coronavirus isolated from a patient with atypical
510 pneumonia after visiting Wuhan. *Emerg Microbes Infect.* 2020;9:221-36.
- 511 39. Xu J, Zhao S, Teng T, Abdalla AE, Zhu W, Xie L, Wang Y, Guo X. Systematic
512 Comparison of Two Animal-to-Human Transmitted Human Coronaviruses: SARS-CoV-2
513 and SARS-CoV. *Viruses.* 2020;12.

514

515 **Tables**

516

	Description	Genbank accession number
Bat-CoV	Bat SARS-like coronavirus isolate As6526	KY417142
Civet-CoV	Civet SARS coronavirus civet007	AY572034
SARS-CoV	Human severe acute respiratory syndrome-related coronavirus strain hTor02	NC_004718
Pangolin-CoV	Pangolin coronavirus isolate PCoV_GX-P2V	MT072864
SARS-CoV-2	Human severe acute respiratory syndrome coronavirus 2 isolate Wuhan-Hu-1	MT072688

517 Table 1. Genome sequences selected for the current study. Note that the SARS-CoV strain
518 hTor02 is from humans infected during the middle and late phases of the SARS epidemic of
519 2013, and has a deletion of 29 nucleotides in the region of ORF8.

520

Name	Bat-CoV			Name	Civet-CoV			Name	SARS-CoV		
	Start	Stop	Length		Start	Stop	Length		Start	Stop	Length
ORF1a*	265	13398	13134	ORF1a	239	13366	13128	ORF1a	265	13398	13134
ORF1b*	13398	21485	8086	ORF1b	13366	21459	8092	ORF1b	13398	21485	8086
S	21492	25217	3726	S	21466	25233	3768	S	21492	25259	3768
ORF3a	25227	26051	825	ORF3a	25242	26066	825	ORF3a	25268	26092	825
ORF3b	25648	25992	345	ORF3b	25663	26127	465	ORF3b	25689	26153	465
E	26076	26306	231	E	26091	26321	231	E	26117	26347	231
M	26357	27022	666	M	26372	27037	666	M	26398	27063	666
ORF6	27033	27224	192	ORF6	27048	27239	192	ORF6	27074	27265	192
ORF7a	27232	27600	369	ORF7a	27247	27615	369	ORF7a	27273	27641	369

ORF7b	27597	27731	135	ORF7b	27612	27746	135	ORF7b	27638	27772	135	
ORF8	27738	28103	366	ORF8	27753	28121	369	ORF8a	27779	27898	120	
								ORF8b	27864	28118	255	
N	28118	29386	1269	N	28123	29391	1269	N	28120	29388	1269	
				ORF9b	28133	28429	297	ORF9b	28130	28426	297	
				ORF9c	28586	28798	213	ORF9c**	28583	28793	211	
	Pangolin-CoV				SARS-CoV-2							
Name	Start	Stop	Length	Name	Start	Stop	Length					
ORF1a	249	13427	13179	ORF1a	251	13453	13203					
ORF1b	13427	21514	8086	ORF1b	13453	21538	8086					
S	21522	25331	3810	S	21521	25369	3849					
ORF3a	25341	26168	828	ORF3a	25378	26205	828					
				ORF3b***	25509	25680	172					
E	26193	26420	228	E	26230	26457	228					
M	26468	27136	669	M	26508	27176	669					
6	27147	27332	186	ORF6	27187	27372	186					
7a	27339	27704	366	ORF7a	27379	27744	366					
7b	27701	27832	132	ORF7b***	27741	27872	130					
8	27839	28202	366	ORF8	27879	28244	366					
N	28218	29471	1254	ORFN	28259	29518	1260					
				ORF9b***	28269	28562	294					
				ORF9c***	28719	28940	222					
				ORF10	29543	29659	117					

521 Table 2. CDS annotations extracted from Genbank, with ORF names standardized according to
522 the SARS-CoV-2 nomenclature.

523 * For convenience, ORF1ab is split into 2 regions corresponding the ORF1ab gene regions
524 upstream and downstream of the frameshift.

525 ** SARS-CoV annotation for ORF9c was propagated from Genbank entry AY274119: SARS-
526 CoV isolate Tor2, where it is annotated as ORF14.

527 *** SARS-CoV-2 annotations for ORF3b, ORF7b, ORF9b and ORF9c were propagated from
528 Genbank entry MN985325: Severe acute respiratory syndrome coronavirus 2 isolate 2019-
529 nCoV/USA-WA1/2020.

530

531

	Predicted: YES	Predicted: NO	Total
Known ORF	11	2	13
Unknown ORF	2	10	12
Total	13	12	25
	Sensitivity=0.85	Specificity=0.83	

532 Table 3. Prediction performance of the GOFIX method on the set of known ORFs in the SARS-
533 CoV genome.

534

Figures

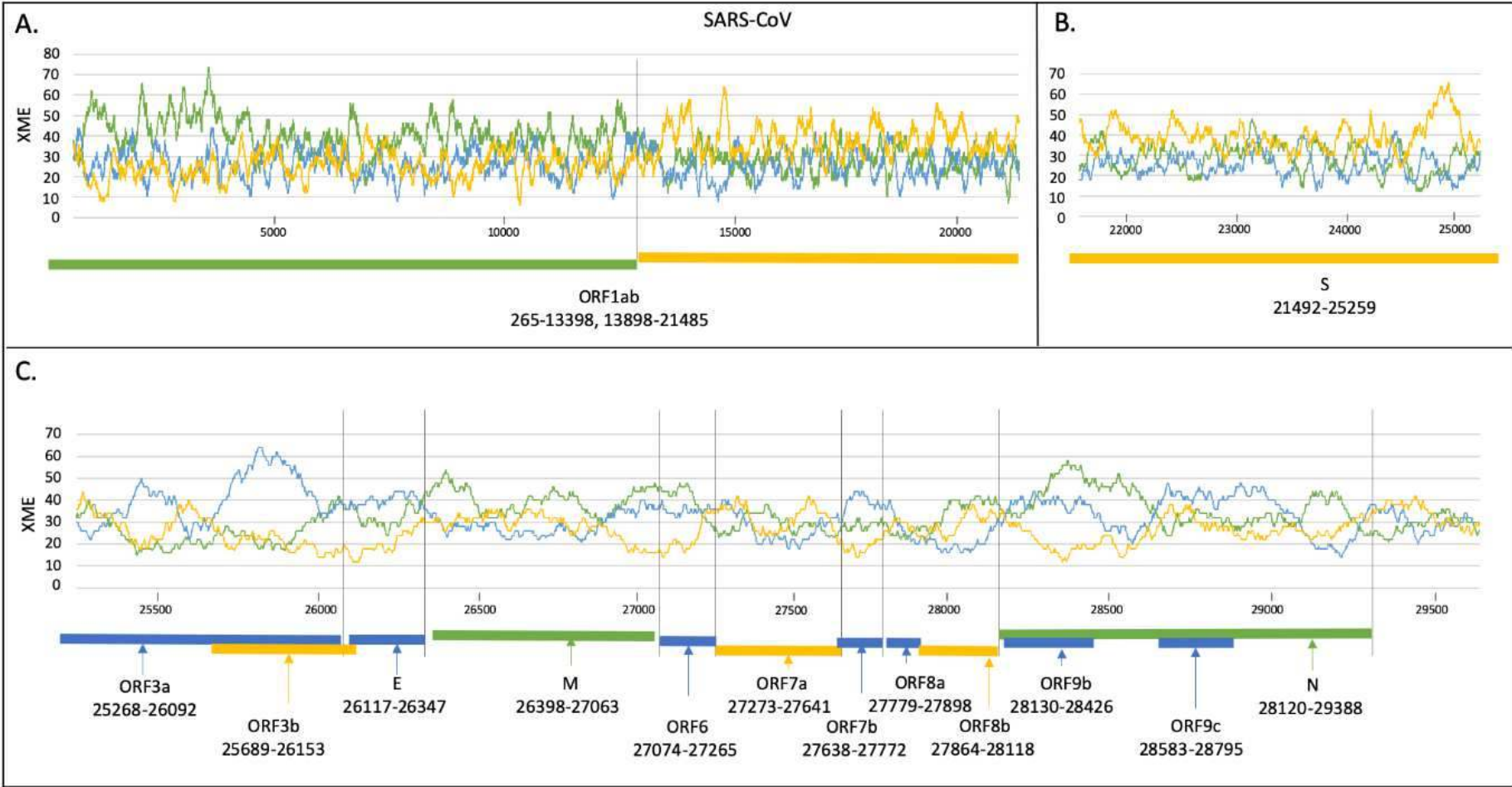


Fig. 1. X motif enrichment (XME_f) scores in the three frames $f=0, 1$ and 2 (green, blue, yellow respectively) of the SARS-CoV genome, using a sliding window of length 150 nucleotides. Genomic organization of known ORFs is shown underneath the plots. **A.** Polyprotein gene ORF1ab. **B.** Spike protein. **C.** C-terminal structural and accessory proteins. The colors used in the enrichment plot and in the boxes representing ORFs (green, blue, yellow) indicate the three frames 0,1 and 2 respectively.

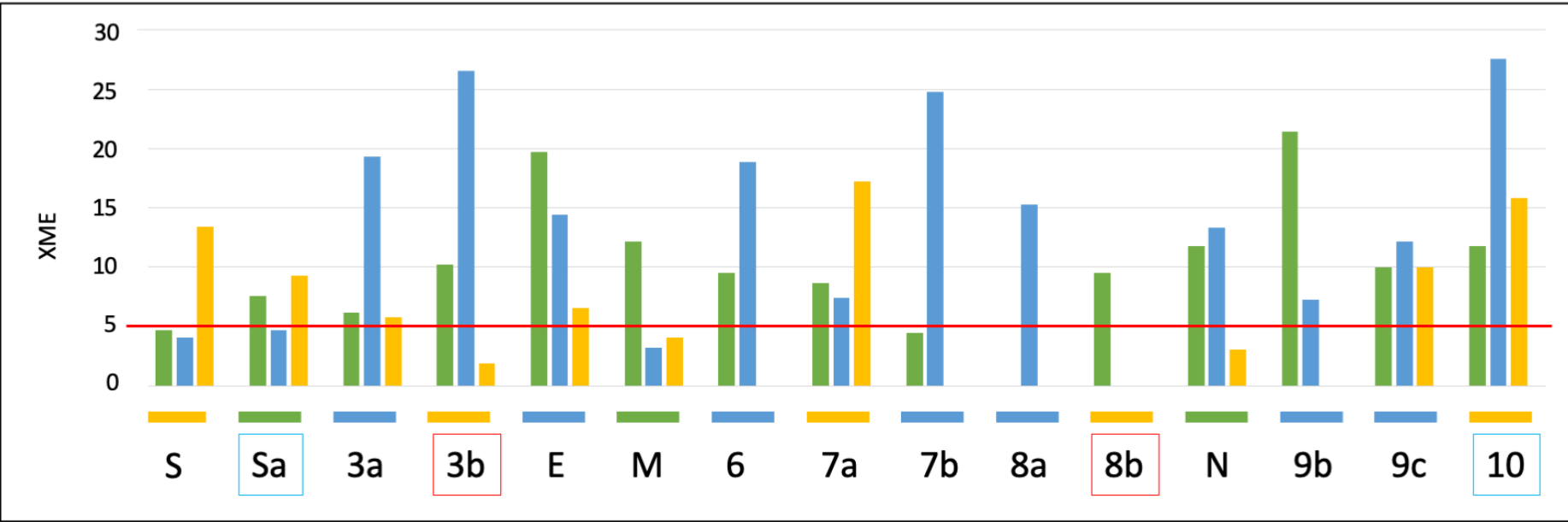


Fig. 2. XME_f scores calculated by GOFIX for potential ORFs in the 3' terminal region of the SARS-CoV genome, in the three frames $f=0, 1$ and 2 (green, blue, yellow respectively). For clarity, only Genbank annotated ORFs or new ORFs predicted in this work are shown. The red line represents the threshold value $XME=XME_f=5$ (where f is the reading frame) for the prediction of a functional ORF. Known ORFs are indicated below the histogram using the color corresponding to the ORF reading frame. Known ORFs not predicted to be functional by GOFIX are outlined in red. Novel ORFs predicted by GOFIX are outlined in blue.

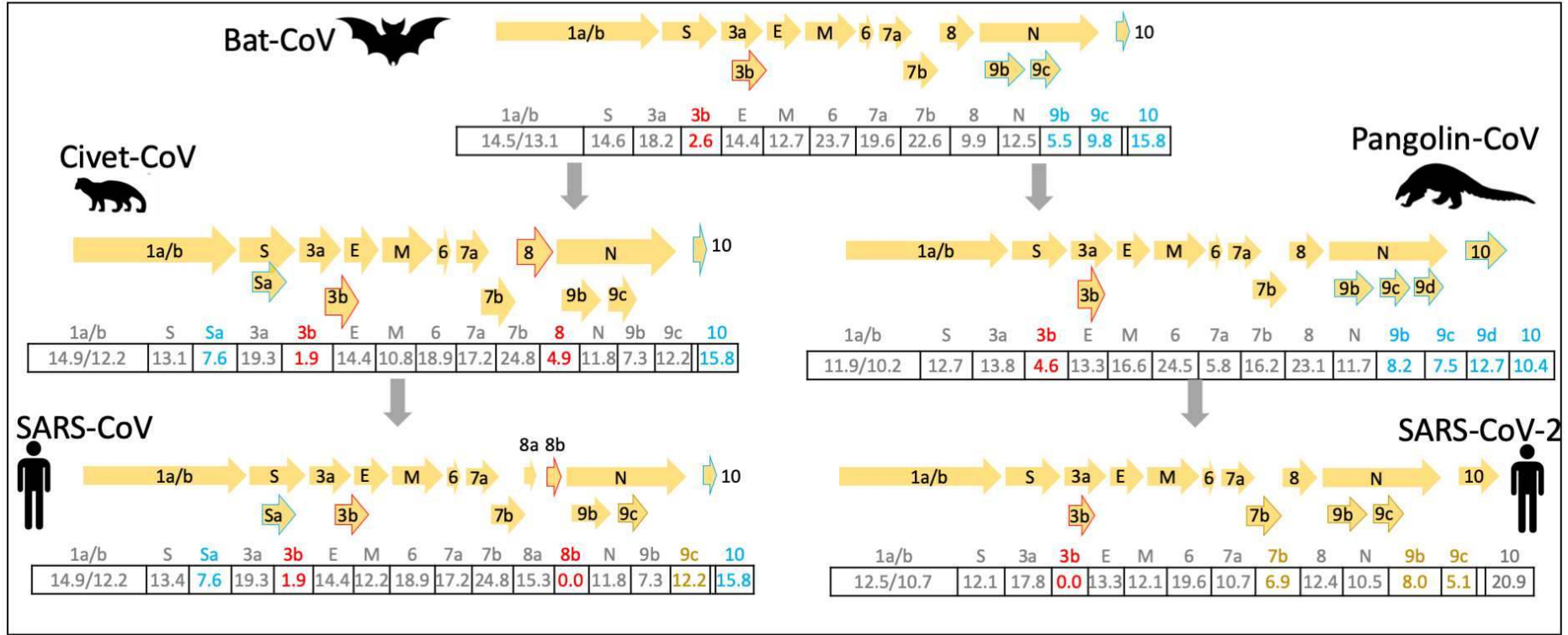


Fig 3. Prediction of ORFs in representative SARS-like coronavirus genomes. A schema is provided for each genome, showing the Genbank annotated ORFs and new ORFs predicted in this work. The numbers in the tables below each schema indicate the XME scores of each ORF. Genbank annotated ORFs that are not predicted to be functional by the GOFIX method are highlighted in red. Novel ORFs predicted by GOFIX are shown in blue. ORFs with conflicting annotations in Genbank, but predicted by GOFIX are shown in brown. Note that ORF3b in Civet-CoV and SARS-CoV is not homologous to ORF3b in Pangolin-CoV and SARS-CoV-2.

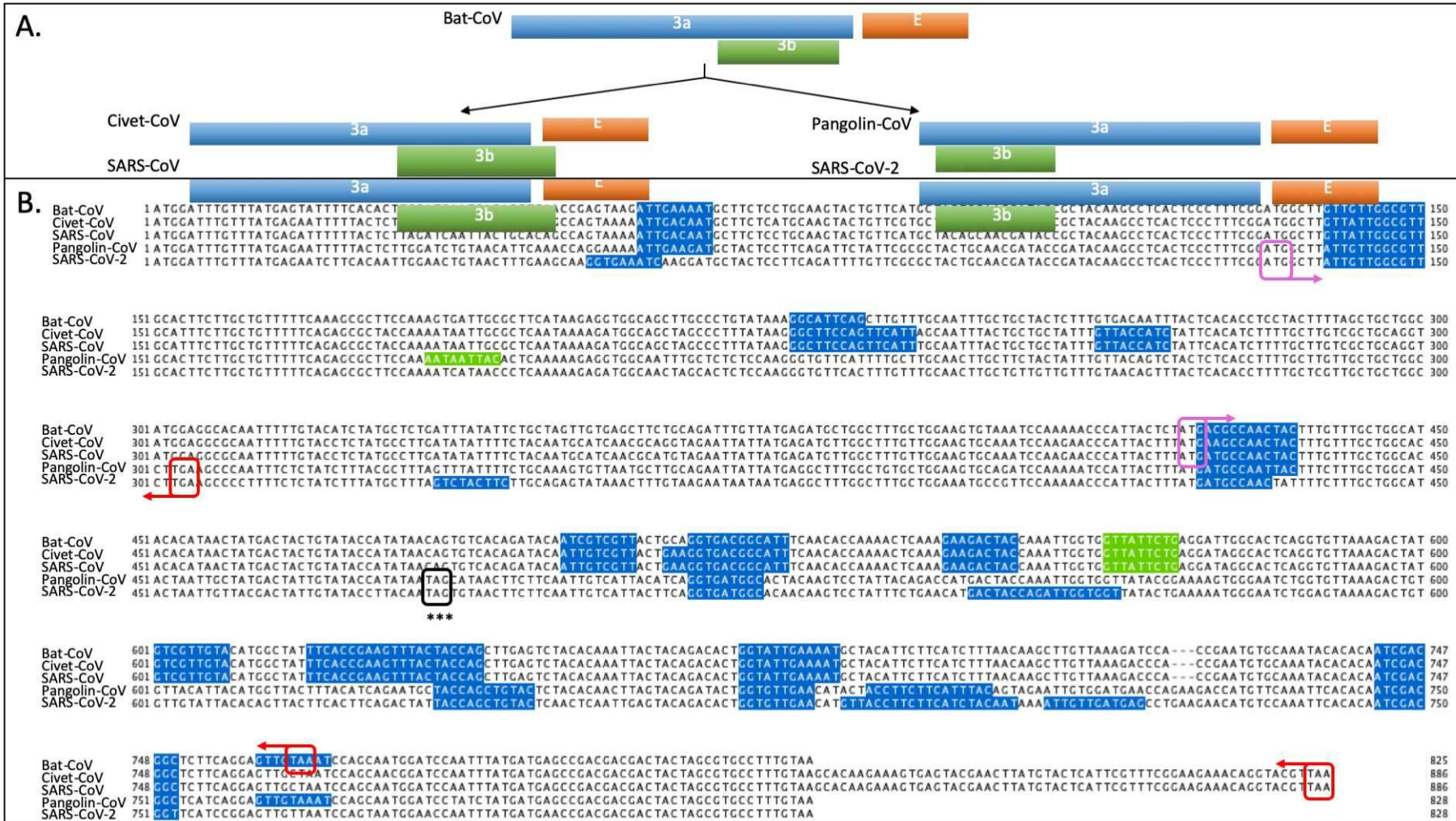


Fig. 4. **A.** Schematic view of genome organization of ORF3a, ORF3b and E gene. **B.** Multiple alignment of ORF3a, ORF3b sequences, with *X* motifs in the reading frame of ORF3a shown in blue. The start and stop codons of the overlapping ORF3b sequences (in the +1 reading frame of ORF3a) are indicated by purple and red boxes respectively. *X* motifs in the reading frame of ORF3b are shown in green.

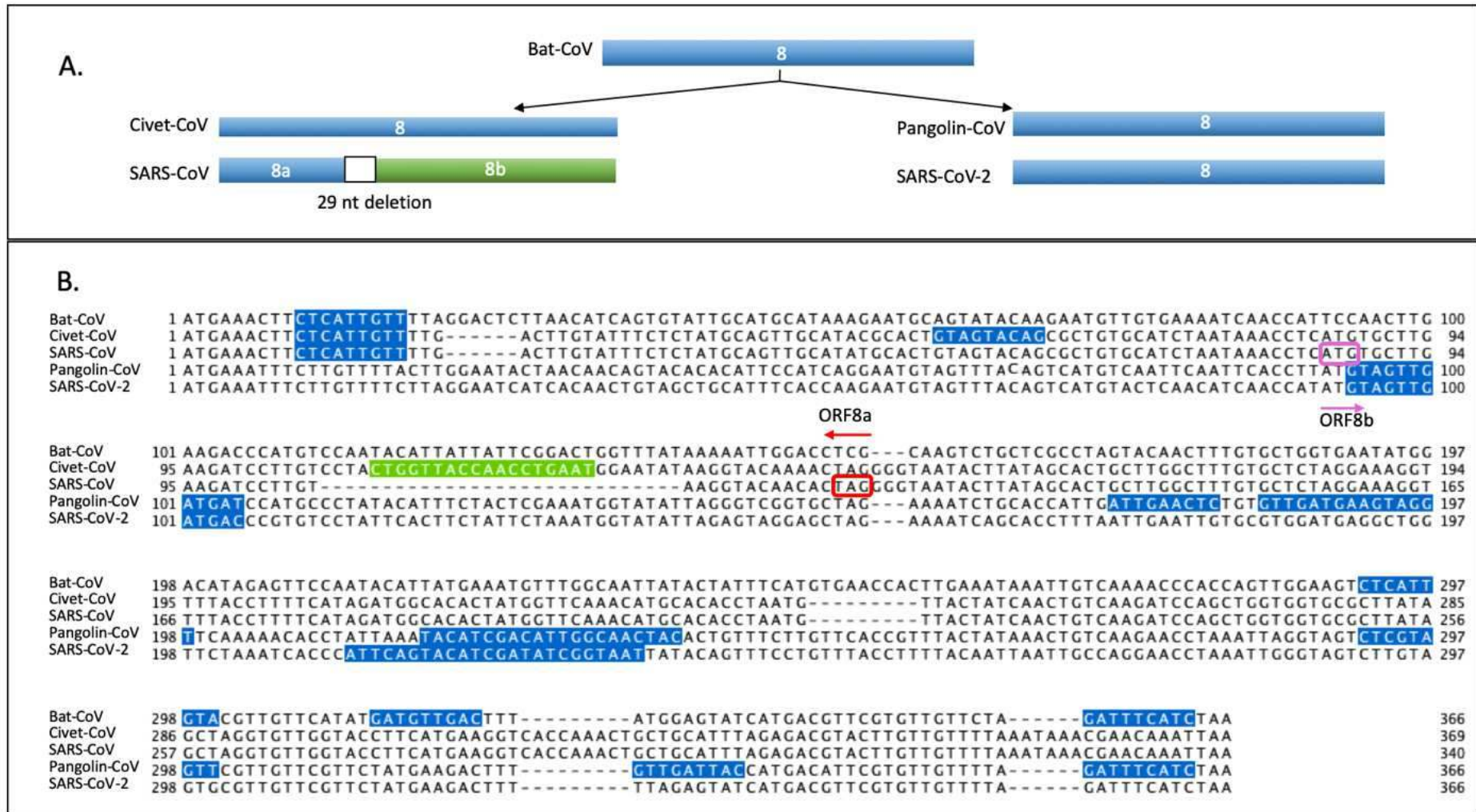


Fig. 5. **A.** Schematic view of genome organization of ORF8, highlighting the 29-nt deletion in SARS-CoV, resulting in 2 ORFs: ORF8a and ORF8b. **B.** Multiple alignment of ORF8 sequences, with X motifs in the reading frame of ORF3a shown in blue. The start and stop codons of the

SARS-CoV ORF8a and ORF8b sequences are indicated by purple and red boxes respectively. The *X* motif corresponding to the 29-nt deletion is shown in green.

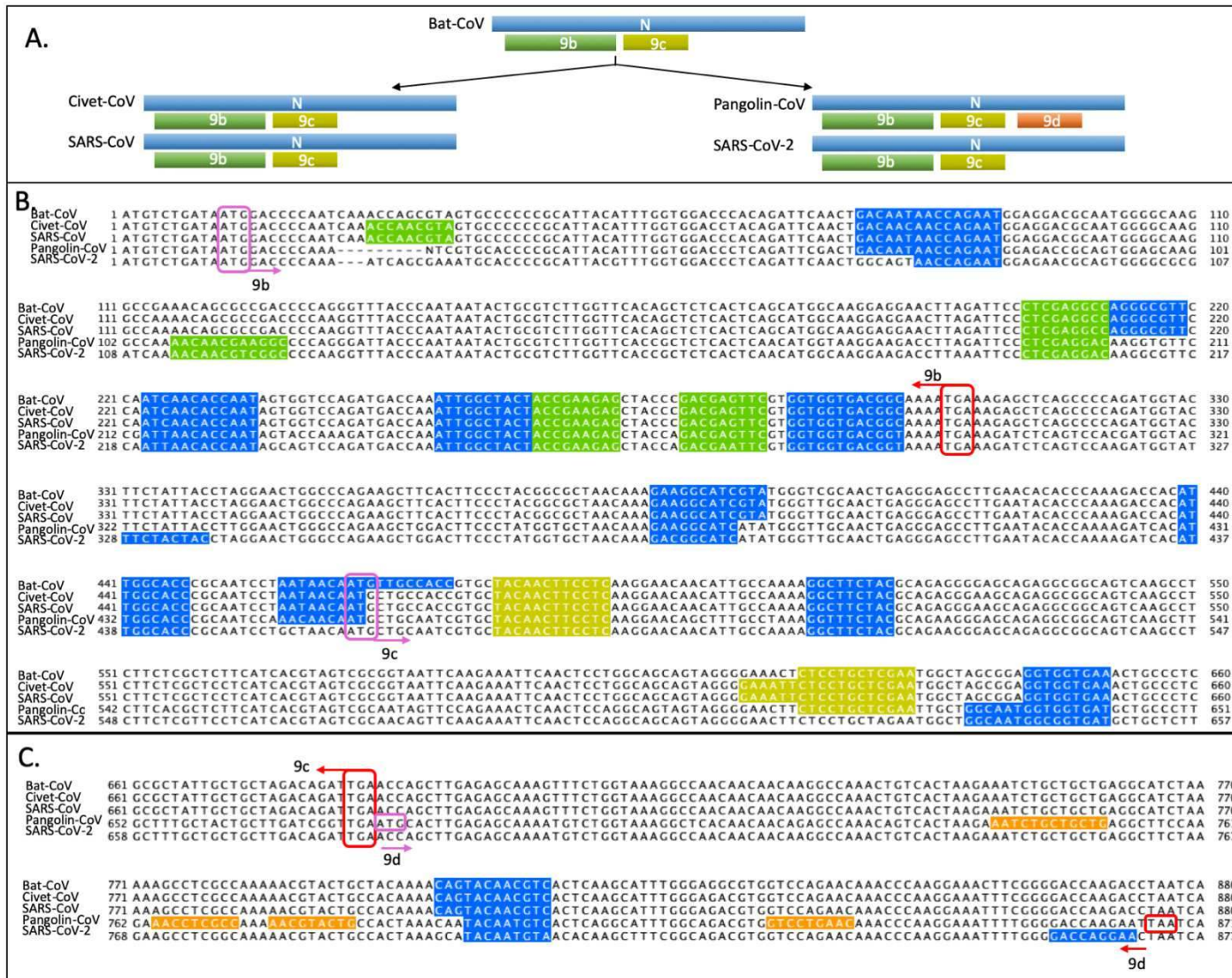


Fig 6. A. Schematic view of genome organization of ORF N, with overlapping genes ORF9b, 9c and the novel predicted 9d. B. Multiple alignment of ORF N sequences, with *X* motifs in the reading frame of ORF N shown in blue, in ORF9b in green, in ORF9c in yellow. Start and stop codons of the overlapping genes are indicated by violet and red boxes, respectively. C. The novel ORF9d predicted in Pangolin-Cov with *X* motifs in the reading frame shown in orange.

Bat-CoV	1 ATGGGCTAT	GTAACGTTTT	GCAATTCGGTTTACGATACATAGTCTACTCTTGTGCAGAATGAATTCTCGTAGCTAAAC	80
Civet-CoV	1 ATGGGCTAT	GTAACGTTTT	GCAATTCGGTTTACGATACATAGTCTACTCTTGTGCAGAATGAATTCTCGTAAC	80
SARS-CoV	1 ATGGGCTAT	GTAACGTTTT	GCAATTCGGTTTACGATACATAGTCTACTCTTGTGCAGAATGAATTCTCGTAAC	80
Pangolin-CoV	1 ATGGGCTAT	GTAACGTTTT	CGCTTTTCCGTTTACGATACATAGTCTACTCTTGTGCAGAATGAATTCTCGTAGCTATAC	80
SARS-CoV-2	1 ATGGGCTATATAA	AACGTTTT	CGCTTTTCCGTTTACGATATATAGTCTACTCTTGTGCAGAATGAATTCTCGTAACTACAT	80
Bat-CoV	81 AGCACAAGTAGGTTTAGTTAACTTTAATCTCACATAG			117
Civet-CoV	81 AGCACAAGTAGGTTTAGTTAACTTTAATCTCACATAG			117
SARS-CoV	81 AGCACAAGTAGGTTTAGTTAACTTTAATCTCACATAG			117
Pangolin-CoV	81 AGCACAAGTAGGTATAGTTAACTTTAATCTCACATAG			117
SARS-CoV-2	81 AGCACAA	GTAGATGAGTTAAC	TTTAATCTCACATAG	117

Fig 7. Multiple alignment of ORF10 sequences, with *X* motifs in the reading frame shown in blue. Stop codons are indicated by red boxes.

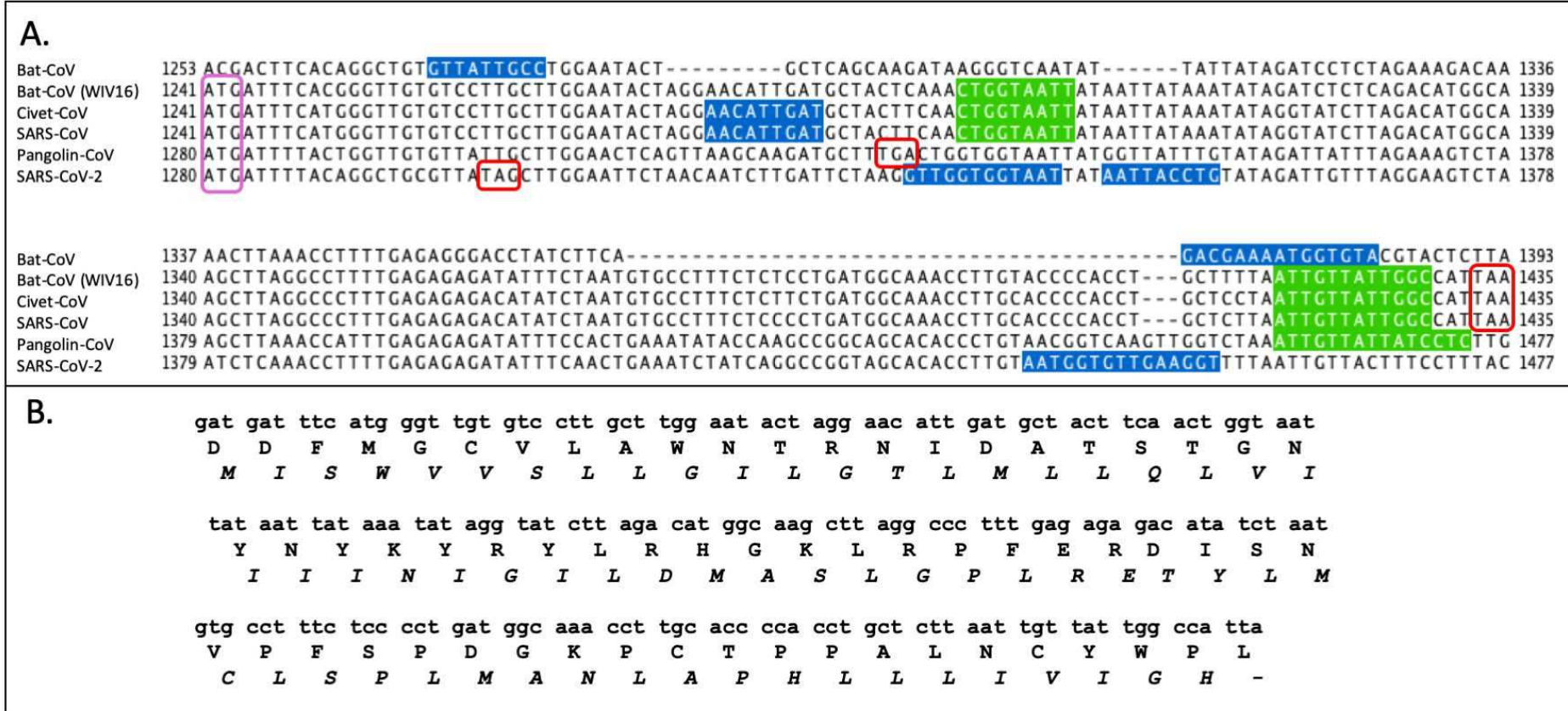


Fig. 8. A. Multiple alignment of ORFSa sequences, with X motifs in the reading frame of ORF S shown in blue and ORFSa in green. Start and stop codons of the overlapping genes are indicated by violet and red boxes, respectively. Bat-CoV (WIV16) sequence is from Genbank:KT444582.

B. Nucleotide and amino acid sequences of the novel ORF predicted to overlap the Spike protein in the genome of SARS-CoV. The nucleotide

sequence segment (SARS-CoV:nt 22732-22926) encodes part (residues 414-478) of the RBD (residues 323-502) of the Spike protein (normal characters), while the reading frame +1 encodes a potential overlapping ORF (*italics*), which we named Sa.

Figures

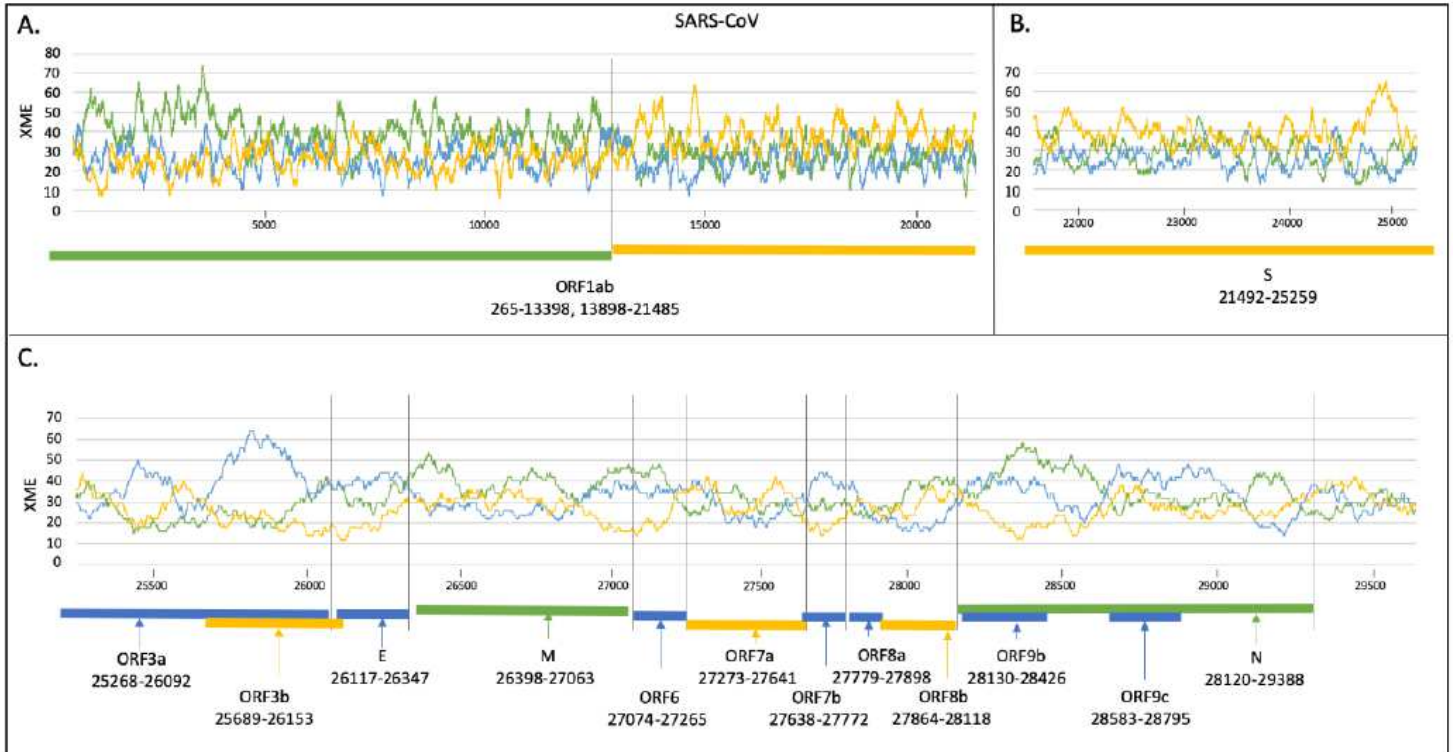


Figure 1

X motif enrichment (XMEf) scores in the three frames $f = 0, 1$ and 2 (green, blue, yellow respectively) of the SARS-CoV genome, using a sliding window of length 150 nucleotides. Genomic organization of known ORFs is shown underneath the plots. A. Polyprotein gene ORF1ab. B. Spike protein. C. C-terminal structural and accessory proteins. The colors used in the enrichment plot and in the boxes representing ORFs (green, blue, yellow) indicate the three frames 0,1 and 2 respectively.

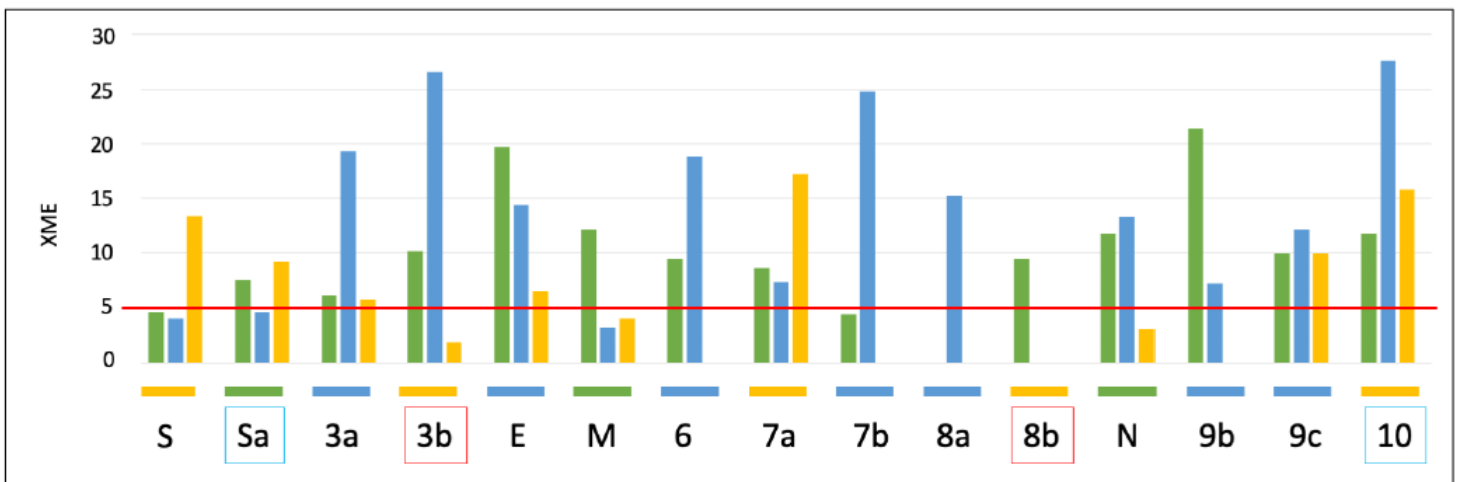


Figure 2

XMEf scores calculated by GOFIX for potential ORFs in the 3' terminal region of the SARS-CoV genome, in the three frames $f = 0, 1$ and 2 (green, blue, yellow respectively). For clarity, only Genbank annotated ORFs or new ORFs predicted in this work are shown. The red line represents the threshold value $XME = XME_f = 5$ (where f is the reading frame) for the prediction of a functional ORF. Known ORFs are indicated below the histogram using the color corresponding to the ORF reading frame. Known ORFs not predicted to be functional by GOFIX are outlined in red. Novel ORFs predicted by GOFIX are outlined in blue.

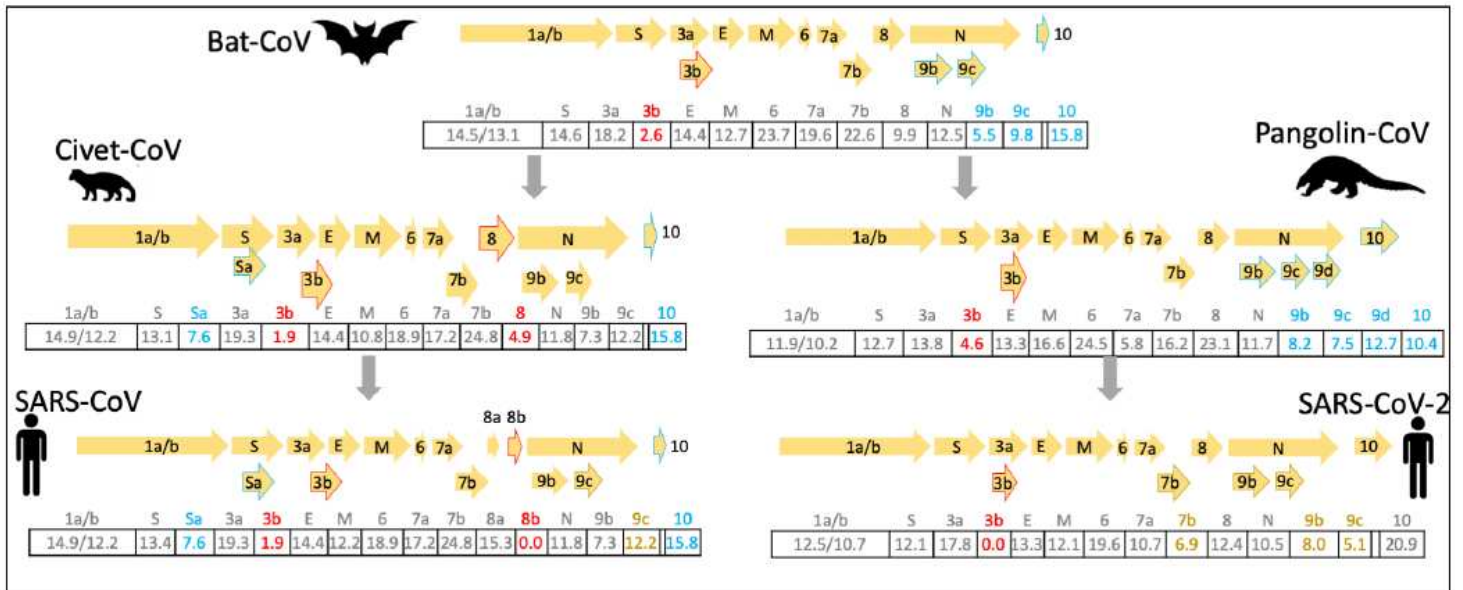


Figure 3

Prediction of ORFs in representative SARS-like coronavirus genomes. A schema is provided for each genome, showing the Genbank annotated ORFs and new ORFs predicted in this work. The numbers in the tables below each schema indicate the XME scores of each ORF. Genbank annotated ORFs that are not predicted to be functional by the GOFIX method are highlighted in red. Novel ORFs predicted by GOFIX are shown in blue. ORFs with conflicting annotations in Genbank, but predicted by GOFIX are shown in brown. Note that ORF3b in Civet-CoV and SARS-CoV is not homologous to ORF3b in Pangolin-CoV and SARS-CoV-2.

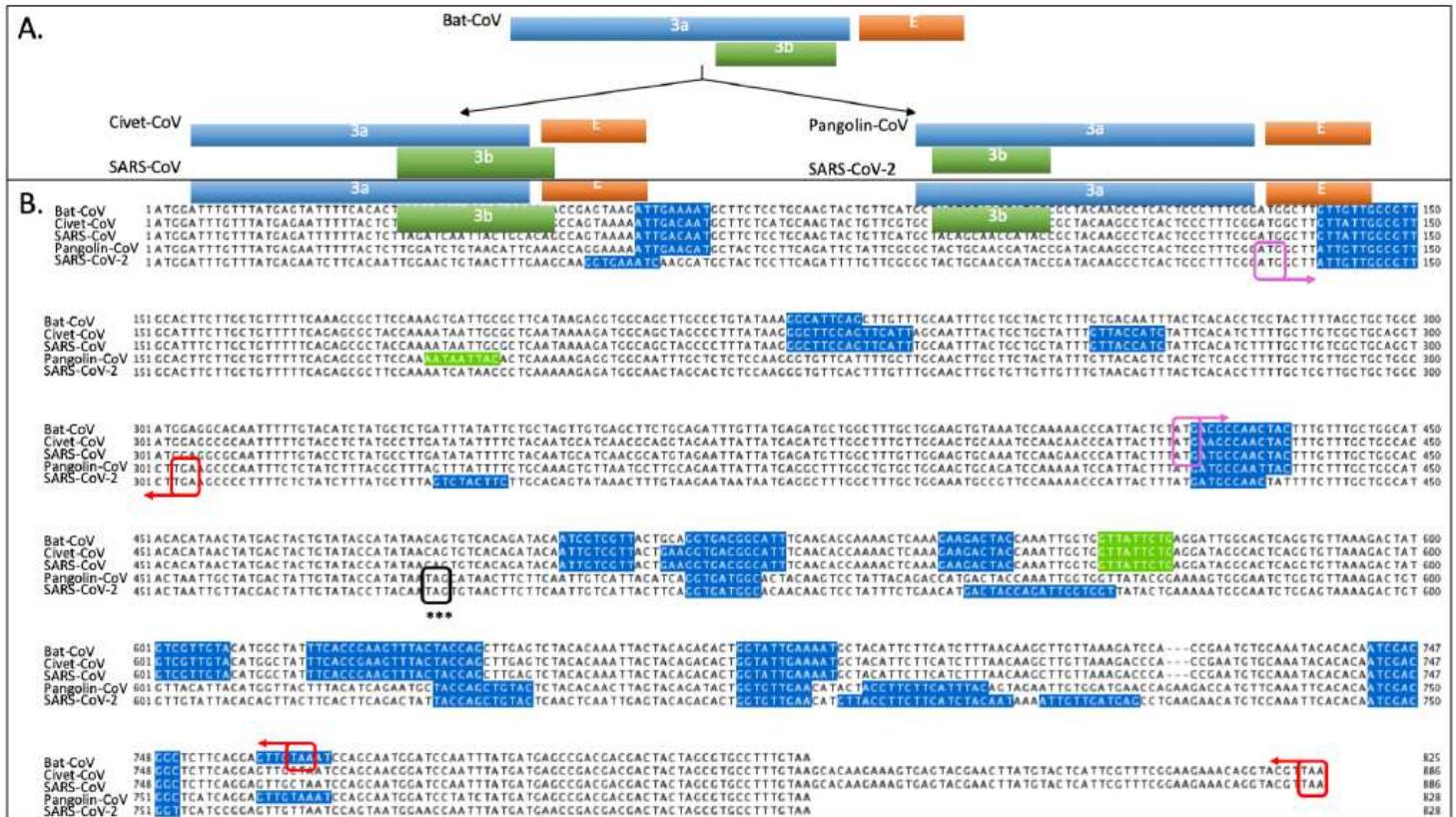


Figure 4

A. Schematic view of genome organization of ORF3a, ORF3b and E gene. B. Multiple alignment of ORF3a, ORF3b sequences, with X motifs in the reading frame of ORF3a shown in blue. The start and stop codons of the overlapping ORF3b sequences (in the +1 reading frame of ORF3a) are indicated by purple and red boxes respectively. X motifs in the reading frame of ORF3b are shown in green.

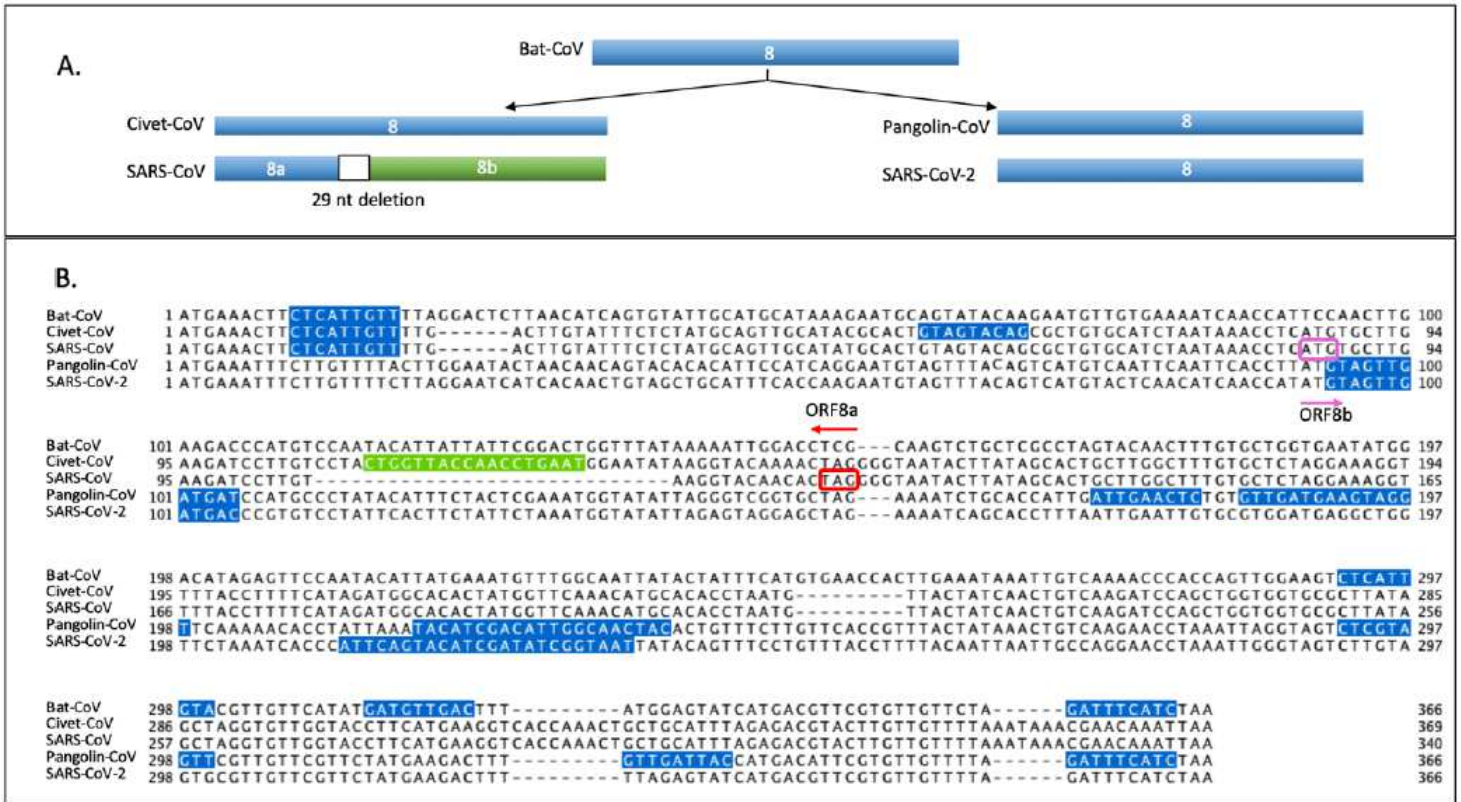


Figure 5

A. Schematic view of genome organization of ORF8, highlighting the 29-nt deletion in SARS-CoV, resulting in 2 ORFs: ORF8a and ORF8b. B. Multiple alignment of ORF8 sequences, with X motifs in the reading frame of ORF3a shown in blue. The start and stop codons of the SARS-CoV ORF8a and ORF8b sequences are indicated by purple and red boxes respectively. The X motif corresponding to the 29-nt deletion is shown in green.

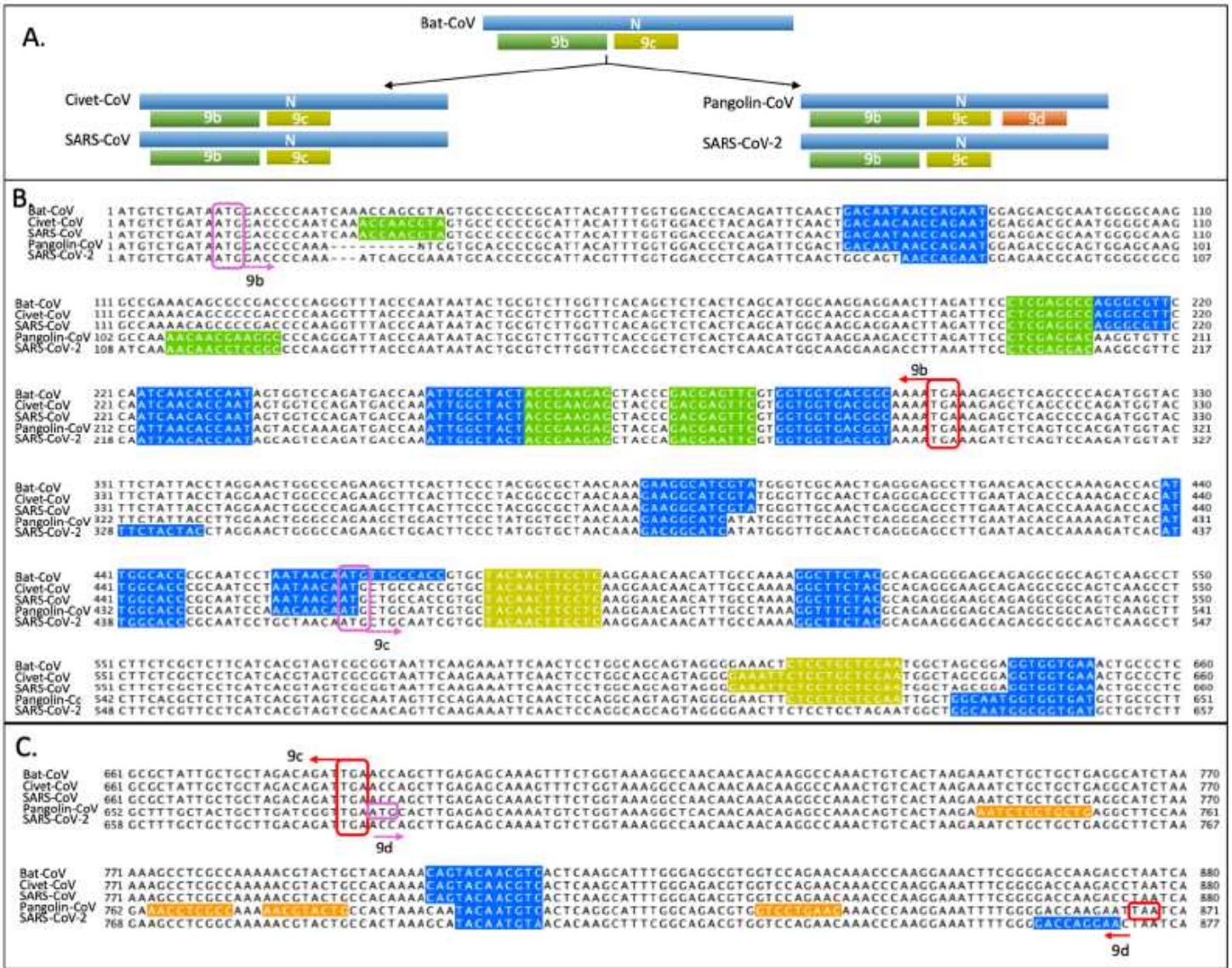


Figure 6

A. Schematic view of genome organization of ORF N, with overlapping genes ORF9b, 9c and the novel predicted 9d. B. Multiple alignment of ORF N sequences, with X motifs in the reading frame of ORF N shown in blue, in ORF9b in green, in ORF9c in yellow. Start and stop codons of the overlapping genes are indicated by violet and red boxes, respectively. C. The novel ORF9d predicted in Pangolin-Cov with X motifs in the reading frame shown in orange.



Figure 7

Multiple alignment of ORF10 sequences, with X motifs in the reading frame shown in blue. Stop codons are indicated by red boxes.

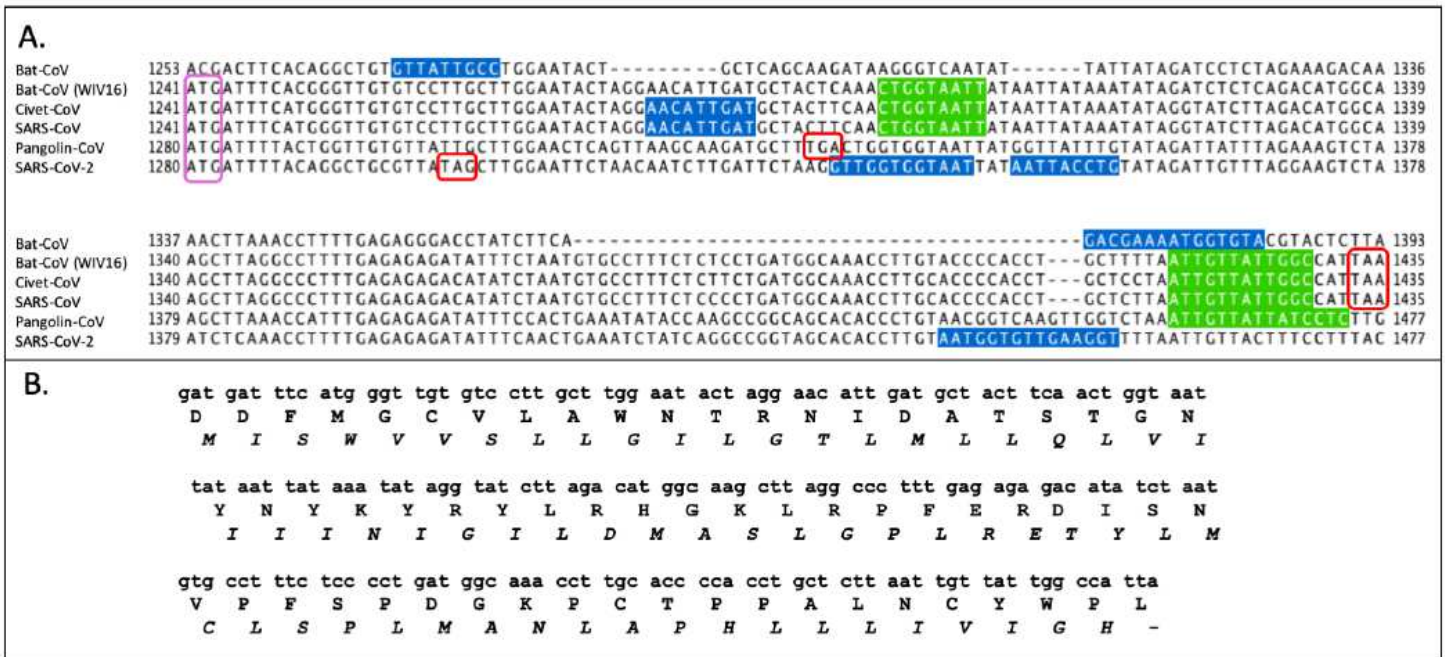


Figure 8

A. Multiple alignment of ORFSa sequences, with X motifs in the reading frame of ORF S shown in blue and ORFSa in green. Start and stop codons of the overlapping genes are indicated by violet and red boxes, respectively. Bat-CoV (WIV16) sequence is from Genbank:KT444582. B. Nucleotide and amino acid sequences of the novel ORF predicted to overlap the Spike protein in the genome of SARS-CoV. The nucleotide sequence segment (SARS-CoV:nt 22732-22926) encodes part (residues 414-478) of the RBD (residues 323-502) of the Spike protein (normal characters), while the reading frame +1 encodes a potential overlapping ORF (italics), which we named Sa.