# Comparisons of the genome of SARS-CoV-2 and those of other betacoronaviruses — Source link ↗

Eduardo Rodríguez-Román, Adrian J. Gibbs

**Institutions:** Venezuelan Institute for Scientific Research, Australian National University

Related papers:

- Mutational analysis and assessment of its impact on proteins of SARS-CoV-2 genomes from India.

- Origin, phylogeny, variability and epitope conservation of SARS-CoV-2 worldwide.

- [The genome comparison of SARS-CoV and other coronaviruses].

- Generation and transmission of interlineage recombinants in the SARS-CoV-2 pandemic.

- Comparison Of SARS-CoV-2 Virus Variant Genomes Detected In China and USA

# Comparisons of the genome of SARS-CoV-2 and those of other betacoronaviruses

Eduardo Rodríguez-Román[1*] and Adrian J. Gibbs[2]

1. *Center for Microbiology and Cell Biology, Instituto Venezolano de Investigaciones Científicas. Caracas 1020A, Venezuela*

2. *Emeritus Faculty, Australian National University, Canberra, ACT 2601, Australia*

**Corresponding author:** Eduardo Rodríguez-Román, PhD

Center for Microbiology and Cell Biology, Instituto Venezolano de Investigaciones Científicas (IVIC). Carretera Panamericana, Km 11. P.O. Box 20632. Caracas 1020A, Venezuela; Tel: +58(212)504-1189; Fax: +58(212)504-1500; Email: erodriguezroman@gmail.com

**Abstract**

The genome of SARS-CoV-2 virus causing the worldwide pandemic of COVID-19 is most closely related to viral metagenomes isolated from bats and, more distantly, pangolins. All are of sarbecoviruses of the genus *Betacoronavirus.* We have unravelled their recombinational and mutational histories. All showed clear evidence of recombination, most events involving the 3' half of the genomes. The 5' region of their genomes was mostly recombinant free, and a phylogeny calculated from this region confirmed that SARS-CoV-2 is closer to RmYN02 than RaTG13, and showed that SARS-CoV-2 diverged from RmYN02 at least 26 years ago, and both diverged from RaTG13 at least 37 years ago; recombinant regions specific to these three viruses provided no additional information as they matched no other Genbank sequences closely. Simple pairwise comparisons of genomes show that there are three regions where most non-synonymous changes probably occurred; the DUF3655 region of the nsp3, the S gene and ORF 8 gene. Differences in the last two of those regions have probably resulted from recombinational changes, however differences in the DUF3655 region may have resulted from selection. A hexamer of the proteins encoded by the nsp3 region may form the molecular pore spanning the double membrane of the coronavirus replication organelle (Wolff et al., 2020), and perhaps the acidic polypeptide encoded by DUF3655 lines it, and presents a novel target for pharmaceutical intervention.

**Keywords:** betacoronaviruses, phylogeny, evolution, SARS-CoV-2, DUF3655, pharmaceutical intervention.

## 1.  Introduction

1

The family *Coronaviridae* is divided into two subfamilies, five genera, 26 subgenera, and 46 species (International Committee on Taxonomy of Viruses; https://talk.ictvonline.org/). However, only members of the genera *Alphacoronavirus* and *Betacoronavirus* have been reported to infect humans.  Coronaviruses (CoVs) have single-stranded positive-sense RNA genomes that are several-fold larger than those of other RNA viruses (Anthony et al., 2017); this reflects the fact that the CoV nsp14 is a proof-reading bi-functional enzyme, ExoN (Ferron et al., 2018) responsible for recombination (Gribble et al., 2020).

In December of 2019 a novel coronavirus causing pneumonia emerged in Wuhan, China (Wu et al., 2020). Initially, the virus was called 2019-nCoV, but it is now known as SARS-CoV-2 (Gorbalenya et al., 2020), and is the etiologic agent of the disease COVID-19.  It is the seventh CoV of humans to be reported (Rodríguez-Román and Gibbs, 2020; Ye et al., 2020), and it has generated a pandemic with more than 10 million people infected, and 0.5 million people dead by the end of June 2020.  While trying to establish from where this virus emerged, there have been conflicting claims that it may have come from bats or pangolins, and is most closely related to either the YN02 virus or the RaTG13 virus (Li et al., 2020; Lin and Chen, 2020; Wang et al., 2020; Xiao et al, 2020; Zhou et al., 2020), and in this short paper we resolve some of these differences and discus an interesting betacoronavirus region - DUF3655!

## Methods

Sequences were downloaded from the Genbank and GISAID databases.  They were edited using BioEdit (Hall, 1999), aligned using the neighbor-joining (NJ) option of ClustalX (Jeanmougin et al., 1998), and the maximum likelihood (ML) method PhyML 3.0 (ML) (Guindon and Gascuel, 2003). Sequences were tested for the presence of phylogenetic anomalies using the full suite of options in RDP4 with default parameters (Maynard-Smith, 1992; Holmes et al., 1999; Padidam et al., 1999; Gibbs et al., 2000; Martin and Rybicki, 2000; McGuire and Wright, 2000; Posada and Crandall, 2001; Martin et al., 2005; Boni et al., 2007; Lemey et al., 2009; Martin et al., 2015); anomalies found by four or fewer methods and with greater than $10^{-5}$ random probability were ignored; statistical support for their topologies was assessed using the SH method (Shimodaira and Hasegawa, 1999). Trees were drawn using Figtree Version 1.3 (http://tree.bio.ed.ac.uk/soft ware/figtree/; 12 May 2018) and a commercial graphics package. Patristic distances within trees were calculated using Patristic 1.0 (Fourment and Gibbs, 2006) to convert trefiles to matrices of pairwise branch lengths.

Pairs of sequences were individually aligned using the TranslatorX server (Abascal et al., 2010; http://translatorx.co.uk).  They were then compared using the DnDscan method (Gibbs et al., 2007), which is a simple heuristic method for scanning aligned sequences, codon-by-codon

38    and codon position-by-position, to identify the NS and S changes that may have occurred
39    converting one codon to the other.  NS and S variation is taken to be the sum of the scores for all
40    pairwise position comparisons within that codon. Each comparison involves substituting a
41    nucleotide of one codon with the homologous nucleotide of the other codon and then checking
42    how this affects the amino acid it encodes using the standard genetic code. The process is then
43    reversed, replacing the nucleotides of the second codon of the pair with those of the first, again
44    only one at a time. Thus, there are six possible exchanges between a single codon pair.  If, say, the
45    first codon is ACT (Thr) and the second GGA (Gly), then all three nucleotides differ and six out
46    of six changed codons are produced.  Substituting of the first position of ACT (Thr) with the first
47    nucleotide of the second codon (GGA) will generate GCT (Ala), a NS change, and similarly
48    substituting of the second position C with the second position G generates AGT (Ser), also a NS
49    change, and  the third generates ACA (Thr), a S change.  Likewise swopping the second GGA
50    (Gly) generates AGA (Arg), GCA (Ala) and GGT (Gly), which are NS, NS and S changes
51    respectively.  In all, the pairwise comparison provides a score of 2/6 S changes, and 4/6 to the NS.
52    Pairs of codons that are identical merely contribute 0/6 to both the total S and NS scores for the
53    window position, and indels are treated as 6/6 NS changes. These calculations make no assumption
54    about the direction of evolutionary change nor of the optimal or most parsimonious path of
55    substitution between two codons. The aim is to assess each of the single possible substitutions
56    indicated by two homologous but different codons. The results for each codon position in the
57    alignment are recorded in a CSV file so that they can be further processed for viewing.  The scores
58    used for Fig. 3, for example, were running (overlapping) sums of 5 codon scores, and thus the
59    NS=5.0 maxima represent five adjacent codons each with a maximum NS score of one.
60         The theoretical isoelectric points of the DUF3655 peptides were calculated using the online
61    ProtParam facility of the ExPASy (Gasteiger et al., 2005; https://web.expasy.org/protoparam/).
62

## Results

64         In mid-May 2020 a BLAST search (Altschul et al., 1990) of the Genbank databases was
65    made using the SARS-CoV-2 Wuhan-Hu-1 sequence (NC_045512) as a query, and over 100
66    related full-length genomic sequences were identified.  These were downloaded, and two from the
67    GISAID database that had been discussed in reports, were added (Rodríguez-Román and Gibbs,
68    2020).

69         The sequences were aligned using MAFFT with its L option.  A Neighbor-Joining (NJ)
70    phylogeny of these sequences identified eight distinct genomic sequences in the SARS-CoV-2
71    lineage, together with eleven others in a more distant divergence that included the SARS-CoV
72    reference sequence (NC_004718), and with an outgroup of ten other coronavirus genomes. These
73    were checked for recombination using the Recombination Detection Program (RDP 4.95) (Martin
74    et al., 2015).  Recombinants were detected in, and between, all betacoronaviruses, but not between
75    them and the outgroup sequences.  Eleven were chosen for analysis; all eight from the SARS-CoV-

76 2 lineage and three from the SARS-CoV's; Table 1 lists their Accession Codes, hosts, source isolate
77 codes, and shortened acronyms, which are used hereafter in this paper and its illustrations.

78       Genes are found in all three reading frames of coronavirus genomes, therefore the 11
79 sequences were aligned using MAFFT-L, and BioEdit (Hall, 1999) was used to create, for each of
80 the eleven, a single concatenated alignment of their open reading frames (i.e. all the genes in the
81 same reading frame). We call these, concats. The concats were aligned, using their encoded amino
82 acids as guide, by the TranslatorX online server (Abascal et al., 2010; http://translatorx.co.uk) with
83 its MAFFT option (Katoh and Standley, 2013), and further refined by hand resulting in a concat
84 alignment of 29,286 nts.

85       The maximum likelihood (ML) phylogeny of the eleven complete sarbecovirus concats
86 (Fig. 1A), calculated by the PhyML method, confirmed that they form two lineages diverging from
87 the midpoint root (circled), one including SARS-1 and the other SARS-2. However, the individual
88 nodes in the SARS-2 crown group were not fully supported statistically in this phylogeny; only an
89 average of 0.89 SH support for the terminal three nodes of the SARS-2 lineage. The concat
90 alignment was therefore checked for recombinants using RDP 4.95, and gave the recombinant map
91 shown in Fig. 2, which shows that all concats have recombinant regions. Notably SARS-2, YN02
92 and RaTG13, which we call the crown group of the SARS-2 lineage, all have two identically placed
93 recombinant regions from the same minor 'parent', Rf4092 (i.e. a SARS-1 lineage bat virus).
94 Significant recombinant regions specific to each these three viruses in the spike region, and 3' to
95 it, provided no additional phylogenetic or dating information as they matched no other sequences
96 in Genbank closely (<84% ID).

97       Concats of the basal branches of the phylogeny, ZC45 and ZXC21, have a large central
98 recombinant region most closely related to the homologous region of HKU3-8, which is of the
99 SARS-1 lineage. Further recombinant regions were found in all the sequences, but mostly in their
100 3' terminal halves and, in summary, only one statistically significant recombinant region (i.e. not
101 marked in Fig. 2 with a black dot at its 5' end) was found between nts 1 and 11496 of all eleven
102 concats, and that was in the ZC45 sequence (nts 1443-1768; parent 'unknown') (Fig. 2). Thus,
103 importantly, the 5' terminal region of all eleven concats, stretching from nts 1 to 11496, was
104 available to obtain a phylogeny based on point mutations alone, and not confounded by
105 recombination; the ZC45 recombinant is unlikely to have distorted the phylogeny much as it is
106 only 2.8% of the 11496 nts.

107       Fig. 1B shows the maximum likelihood (ML) phylogeny calculated from nts 1-11496
108 region of the 11 concats. All nodes in this phylogeny have full statistical support (i.e. 1.0 SH), and
109 most of the 'root to tip' distances in the tree were similar, unlike those in Fig. 1A; one effect of
110 recombination. The topology of the '1-11496' tree was different from that of the tree of complete
111 concats as the SARS-2 concat now groups with all SARS-2 lineage bat isolates, and the pangolin
112 isolates are now basal. Closest to the SARS-2 concat is the YN02 concat with the RaTG13 concat
113 a little further away.

114 The minor 'parent' of the shared recombinant regions in the centre of the SARS-2, YN02
115 and RaTG13 concats (nts 14372-15124 and nts 16383-17566) is Rf4092 of the SARS-1 lineage.
116 These recombinant regions were not found in the other bat sequences of the SARS-2 lineage
117 indicating that they resulted from a recombination event that occurred after the crown group
118 diverged from the ZXC21 and ZC45 branch, but before RaTG13 diverged.  Confusingly however
119 the second of these recombinant regions was also found in the pangolin G/1/19 concat!

120 The recombination map also shows the complex recombinational history of the spike gene,
121 the position of which is coloured yellow in the simplified genomic map at the top of Fig. 2.  This
122 is confirmed by the phylogeny of that region (Fig. 1C), which is fully supported statistically except
123 for the SARS-1, Rf4092 and HKU3-8 cluster (mean 0.91 SH).  The spike phylogeny has pangolin
124 genes immediately basal to the SARS-2 and RaTG13 twig, and the spike region of YN02 gene is
125 shown to be from the SARS-1 lineage.  However, it is essential to realize that, although we know
126 the hosts from which the isolates were collected, other hosts may have been infected *en route*.

127 The dates of the nodes in the '1-11496' SARS-2 phylogeny (Fig. 1B) can be inferred using
128 published estimates of the evolutionary rate of the SARS-2 population in the human population,
129 assuming that the pre- and post- emergence rates are the same (Rodríguez-Román and Gibbs,
130 2020).  Various estimates of the SARS-2 evolutionary rate have been published recently; 1.126 x
131 $10^{-3}$ (95 % BCI: 1.03–1.23 x $10^{-3}$) substitutions per site per year (s/s/y) (Candido et al., 2020),
132 $1.1 \times 10^{-3}$ s/s/y (95% CI $7.03 \times 10^{-4}$ and $1.5 \times 10^{-3}$ s/s/y) (Duchêne et al., 2020), $9.41 \times 10^{-4}$ s/s/y +/-
133 $4.99 \times 10^{-5}$ (Pybus et al., 2020) and 8 x $10^{-4}$ s/s/y (Resende et al., 2020).

134 The mean of these rate estimates is 0.99 x$10^{-3}$ s/s/y, and, assuming that the virus is
135 evolving at the same rate as the '1-11496' region of its genome, then the mean patristic distances
136 passing through nodes in Fig. 1B suggest that the SARS-2 and YN02 viruses diverged in 1994
137 CE (26.03 years before present; ybp), they diverged from RatG13 in 1983 CE (36.8 ybp), and
138 from ZXC21 and ZC45 in 1936 CE (83.6 ybp) and from G/1/19 in 1908 CE (111.8 ybp). The
139 standard deviation of the branch length estimates varied between 0.8% and 2.8%.  The most
140 recent estimates are probably the most accurate because although all mutations contribute to the
141 'molecular clock', most are quickly lost (Duchêne et al. 2014), and therefore times to the older
142 dates are overestimated. Nonetheless, it is probable that SARS-2 and YN02 diverged over 20
143 years ago, and the two recombinant regions characteristic of the SARS-2, YN02 and RaTG13
144 virus genomes were acquired by their shared progenitor more than 30, but less than 80, years
145 ago!  All these datings are based on a large number of assumptions, and could be earlier as
146 concluded by Wang et al. (2020).

147 Finally, we compared the concat sequences directly in pairs, not only to identify any
148 regions that were evolving abnormally, but also to confirm the recombination map patterns shown
149 in Fig. 2.  We used the DnDscan method (Gibbs et al., 2007 - see Methods) as this enables simple
150 visual comparisons to be made, as well as numerical.  Fig. 3 shows the synonymous (S - blue) and
151 non-synonymous (NS - gold) differences in five of 45 pairwise possible comparisons of eleven

152    concats. It can be seen that S differences occur throughout most of the comparisons, but NS
153    differences are most obvious in three regions of the genomes. There are slightly fewer S
154    differences between SARS-2 and YN02 than between SARS-2 and RatG13 or between YN02 and
155    RaTG13, and this confirms the phylogenetic tree (Fig. 1B); it shows that SARS-2 is closest to
156    YN02 as the total DnDscan scores for the '1-11496' regions of SARS-2 v YN02 are S 100.0 NS
157    27.5, but, for the other combinations, S134.0 NS 30.1 and S131.5 NS37.6, respectively. The
158    largest NS differences are in the DnDscans of the spike protein gene, especially its RBD region
159    and an adjacent "-PRRA- " insertion (Andersen et al. 2020). Again, the recombination map results
160    (Fig. 2) are confirmed by the DnDscan as the spike region of the SARS-2 x RaTG13 comparison,
161    especially its 5' end, has few NS differences, as they share an 'unknown' recombinant (nts 21257 -
162    22152).

163        There are also two other regions of the concats consistently showing larger numbers of NS
164    differences. One is centred on the 'Domain of Unknown Function' (DUF) 3655 region of the nsp3,
165    a "disordered binding region" (Prates et al. 2020), that is N' terminally adjacent to the ADP-ribose
166    phosphatase. This region of increased NS differences was found to some extent in all concat
167    comparisons suggesting that its differences result from evolution/selection, whereas the other, the
168    ORF 8 region near the 3' end of the genome, was not found in some comparisons, such as SARS-
169    2 x RaTG13, and may therefore have resulted from recombination. The DUF3655 region is
170    discussed below.

171

## 172    2. Discussion

173    We have discombobulated the recombinational and mutational history of the SARS-2 lineage
174    of betacoronaviruses and their metagenomes using the published genomic sequences, despite the
175    possibilities, in this metagenomic age, of the sort of problems outlined by Chan and Zhan (2020).
176    We have shown that the 5' third of their genome is largely free of recombinant regions, n-rec,
177    whereas the remainder is a mélange of recombinant regions from various 'parental' genomes. The
178    SARS-2 crown group share a distinctive pair of recombinant regions that are most closely related
179    to the homologous region of the SARS-1 lineage bat virus, Rf4092. A phylogeny calculated from
180    the 5' n-rec region of the eleven concats shows that the SARS-2 lineage has basal branches of
181    viruses isolated from pangolins, and a crown group consisting of SARS-2 together with YN02,
182    RatG13, ZXC21 and ZC45 all of which come from bats, and in that phylogeny SARS-2 is more
183    closely related to YN02 than RaTG13. This is confirmed by the DnDscan comparisons of the
184    three viruses. YN02, however, has a recombinant region in its 3' half (nts 21098-24042) of
185    'unknown' parentage, but which is probably close to the pangolin virus G/1/19, and which is not
186    present in SARS-2 or RaTG13. Thus, most of the SARS-2 concat, especially its 5' 39%, is closest
187    to the homologous regions of YN02, but the intact concats of SARS-2 and RaTG13 are more
188    distant but complete.

189        Our conclusions about the relationships of the SARS-2 crown group are confirmed in the
190    report of Latinne et al (2020; Fig. 3A) of a large survey of bat viruses of SE Asia. They used
191    primers to amplify a 440 nts region of the RdRp genes of these viruses, and based their phylogeny
192    on that region. Although their amplicon overlapped the 3' end of one of the recombinant regions
193    shared by the SARS-2, YN02 and RaTG13 concats, the overlap is only 78 nts (18%), and the
194    comparison of the 440 nts amplicons found SARS-2 to be closest to YN02.

195        DnDscan, a simple direct comparison of two sequences, found regions of NS change where
196    other more complex methods (Angeletti et al., 2020) did not, and we overcame possible problems
197    with sliding windows (Schmid and Yang, 2008) by making several homologous comparisons. The
198    NS differences around codon 1000 of the DnDscans are from the DUF3655 region. DUF3655
199    marks the 5' end of the nsp3 region and is adjacent to its ADP-ribose phosphatase gene (Michalska
200    et al., 2020). It encodes the N-terminal portion of the nsp3 protein, which has recently been
201    identified by cryo-electron microscopy as forming hexameric molecular pores spanning the double
202    membrane of the coronavirus replication organelle (Wolff et al., 2020). The pores probably allow
203    the progeny SARS-2 genomes to pass from the replication organelle into the lumen of the cytosol,
204    where their 'structural genes' are translated, and together they are assembled to form progeny
205    virions (Hsin et al., 2018). Table 2 shows the DUF3655 peptides of the eleven betacoronaviruses
206    with the acidic and basic residues outlined with different colours; acidic residues in red, and the
207    few basic residues in blue, and with the theoretical pI of these peptides ranging from 3.01 - 3.40,
208    in sharp contrast to the nucleocapsid protein encoded by ORF9 which binds the progeny genomes
209    in the cytosol and has a pI of 10.07 (McBride et al., 2014; Verheije et al., 2010). Table 2 also
210    shows the secondary structures of the SARS-2 crown group DUF3655 proteins predicted by the
211    PSIPRED Workbench (Buchan and Jones, 2019). The DUF3655 proteins are found to have similar
212    N-terminal regions of unstructured residues attached to homologous helical regions, and with C-
213    termini that are more variable in length and composition. The fact that the DUF3655 protein is so
214    acidic indicates its likely function in the pore where it may both electrostatically stabilize the lumen
215    of the pore (Desikan et al., 2020) and ensure that long negatively charged nucleic acid molecules,
216    like progeny viral genomes, are held centrally in the lumen of the pore as they pass through.

217        The FFPred Prediction database of PSIPRED (Cozzetto et al., 2016) found that the most
218    likely "biological process" of SARS-2, YN02 and RaTG13 that involves their DUF3655 proteins
219    is "regulation of metabolic process" (mean probability 0.978) and "regulation of gene expression"
220    (0.908), their "molecular function" is "nucleic acid binding" (0.966) and "DNA binding" (0.890)
221    and their "cellular compartment" is "membrane" (0.785).

222        The DUF3655 region seems to have evaded virological, medical and pharmaceutical
223    scrutiny so far (e.g. Chen and Zhong, 2020; Wei et al., 2020). We suggest that it is probably
224    involved in a unique rate-limiting step of the coronavirus replicative cycle, and may make CoV
225    infections susceptible to drugs, like chloroquine, that increase cellular pH
226    (https://www.sciencemediacentre.org/expert-reaction-to-questions-around-potential-treatments-
227    for-covid-19/ March 18 2020). The detailed analysis of this region, specially from residues 9 to

228    27, which have many negatively charged amino acids (Asp and Glu) (Table 2), and probably the
229    absence of binding sites for macromolecules (RNA, DNA and proteins), would suggest that this
230    region might be an excellent target for the development of an effective treatment for
231    sarbecoviruses.

232        The DUF3655 region warrants more attention especially as repetitive acidic amino acids
233    are present in similar regions of the genomes of human αCoV (JX504050, KF514433, MT438700),
234    MERS-βCoV (MN481964), bulbul δCoV (NC_011547) and infectious bronchitis γCoV
235    (NC_001451).

236

237

238    **Legends**

239    **Fig. 1.** Maximum likelihood phylogenies of eleven sarbecoviruses calculated from A) their
240    complete concat sequences; B) only nts 1-11496 of the concat (i.e. the recombinant-free 5' end);
241    C) the spike protein genes (nts 21315-25143). Acronyms as in Table 1, human viruses in red, bat
242    viruses in blue and pangolin viruses in gold. Midpoint root circled. All nodes have 1.0 SH support
243    except, in Fig. 1A, the three terminal nodes of the SARS-2 lineage (mean 0.89 SH) and, in Fig 1C,
244    the terminal node of the SARS-1 lineage (0.84 SH).

245    **Fig. 2.** Screenshot of the recombinant map of eleven betacoronaviruses analysed using the RDP
246    version 4.95 program with, above, a simplified genome map showing the positions (yellow) of the
247    DUF3655, spike and ORF8 genes.  The recombinant segments that are statistically supported by
248    fewer than five methods and $e^{-5}$ mean probability have a black circle at their 5' end.

249    **Fig. 3.** DnDscan histograms of five pairs of complete betacoronavirus concats; each bar is the
250    running sum of five S (blue) and NS (gold) codon scores with, above, a simplified genome map
251    showing the positions (yellow) of the DUF3655, spike and ORF8 genes.

252

253 **Table 1.** Sarbecovirus genomes compared in this study

| Accession Code | Host | Isolate (acronym) | Country |
|---|---|---|---|
| EPI_ISL_410721 | Pangolin | Guangdong/1/2019 (G/1/19) | China |
| EPI_ISL_412977 | Bat | RmYN02 (YN02) | China |
| GQ153543 | Bat | HKU3-8 (HKU3-8) | HK |
| KY417145 | Bat | Rf4092 (Rf4092) | China |
| MG772933 | Bat | ZC45 (ZC45) | China |
| MG772934 | Bat | ZXC21 (ZXC21) | China |
| MN996532 | Bat | RaTG13 (RaTG13) | China |
| MT040333 | Pangolin | GX-P4L (GX-P4L) | China |
| MT040336 | Pangolin | GXP5E (GXP5E) | China |
| NC_004718 | human | Tor2 (SARS-CoV) (SARS-1) | Canada |
| NC_045512 | human | Wuhan-Hu-1 (SARS-2) | China |

254

255

**Table 2.** Comparison of the DUF3655 region of the eleven betacoronaviruses analysed in this study

```
Sequence                                                                                                          pI

SARS-2   MYCSFYPPDEDEEEGDCEEEFFPSTQY--FYGTFDDYQGKPLFFGATS-AAL-QPEEEQEEDWLDDDSQQTVGQQRGSFDNQTTTIQTIVFVQPQLFMFLTPVVQ-TIF-VN 3.03

YN02     MYCSFYPPDEDEEEGECEEEFFPSTQY--FYGTFDDYRGKSLFFGATS-AAP-QPEEEQEEDWLDDASQQTVAQF-MSGLNQTT-IQSIVFVQPQLFMFPTPVVQ-TIF-VN 3.26

RaTG13   MYCSFYPPDEDEEEGDCFEEDFFPPTQY--FYGTFDDYQGKSLFFGATS-VTP-QPFEELEEDWLDDDSQQTVVQEDDSFVNQTTITQSIAFVQPQLFMFPTPVVQ--TF-VN 3.01

ZC45     MYCSFYPP-FDEGEDDCEFGQCFPSTQY--FYGTFDDYQGKPLFFGATSFSSS-SQFEEQEEDWLFSFSQD--GQFTAV-FNKI----SSVFVPPVLQVFFSTPVVTFTSF-QN 3.31

ZXC21    MYCSFYPP-EDEGEDDCEEGQFFPSTQY--FYGTFDDYQGKPLFFGATSFSSS-SQEEEQEEDWLFSFSQD--GQFT------------------------AVTKTSF-QN 3.36

G/1/19   MYCSFYPPDEDYEEDECEEEQYFPSTQY--FYGTFDDYQGKSLFFGSTS-SAS-QIEEEPEEDWLELGNEEIAMQF------QT----STVFVQSQ-EILSTPVVSEINFSVN 3.08

GX-P4L   MYCSFYPPDEDYEEFYSEEFQPFQPTQY--FYGTFSDYKGLPLFFGASS-V---QQQFFQEEDWLFTFAFV-VEQFVTPTFQEEFL--SITFIVP--AVFQTTIVF--LF-CD 3.11

GX-P5E   MYCSFYPPDEDYEEFYSEEFQPFQPTQY--FYGTFSDYKGLPLFFGASS-V---QQQFFQEEDWLFTFAFV-VEQFVTPTFQEEFL--SITFIVP--AVFQTTIVF--LF-CD 3.11

HKU3-8   MYCSFYPPDEEEDCEECEDEEFISEETCFHFYGTFDDYKGLPLFFGAST-ETPHVEEEEEEELWLDLAIEA----FSFP-------------------FPLP-----EEPVN 3.37

Rf4092   MYCSFYPPDEEEDCDFYDEEEFVPEFSCAHFYGTFEDYRGLPLFFGAST-FM--QVEEEEEEDWLGFATFL-SFHFLFP-------------------FLTP-----EEPVN 3.40

SARS-1   MYCSFYPPDEEEEDDAFCEEEFI-DFTCFHFYGTFDDYQGLPLFFGASA-FTVR-VFEEEEEDWLDFTTFQ-SFIFPFP-------------------------FPTPFFPVN 3.21
```
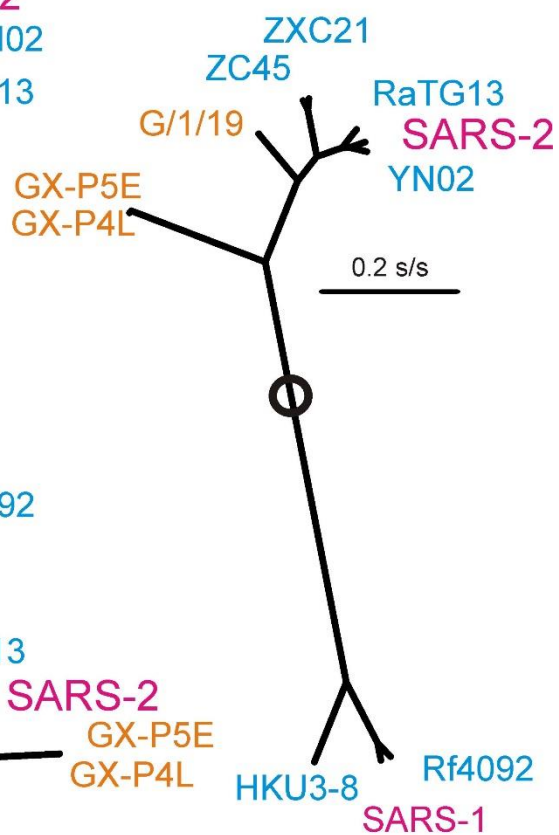
Secondary structure

```
SARS-2   CCCCCCCCCCCCCCCCCCCCCCCCCCCEE--EECCCCCCCCCCCCCCCCC-CCC-CCHHHHHCCCCCCHHHCCCCCCCCCCCCCCCEEEEEEEEECCCEEEEEECCEE-EEE-CC

YN02     CCCCCCCCCCCCCCCCCCCCCCCCCCCEE--EECCCCCCCCCEECCCCCC-CCC-CCHHHHHHHHHHHHHHHHHHHH-CCCCCCEE-EEEEEEECCCCCCCCCCCCEE-EEE-EC

RaTG13   CCCCCCCCCCCCCCCCCCCCCCCCCCCE--ECCCCCCCCCCCEEEECCEE-CCC-CCHHHHHHHCCCCHHCEEEECCCCHHHHHHHEEEHHHCCCCCCCCCCCCCEE--EE-CC
```
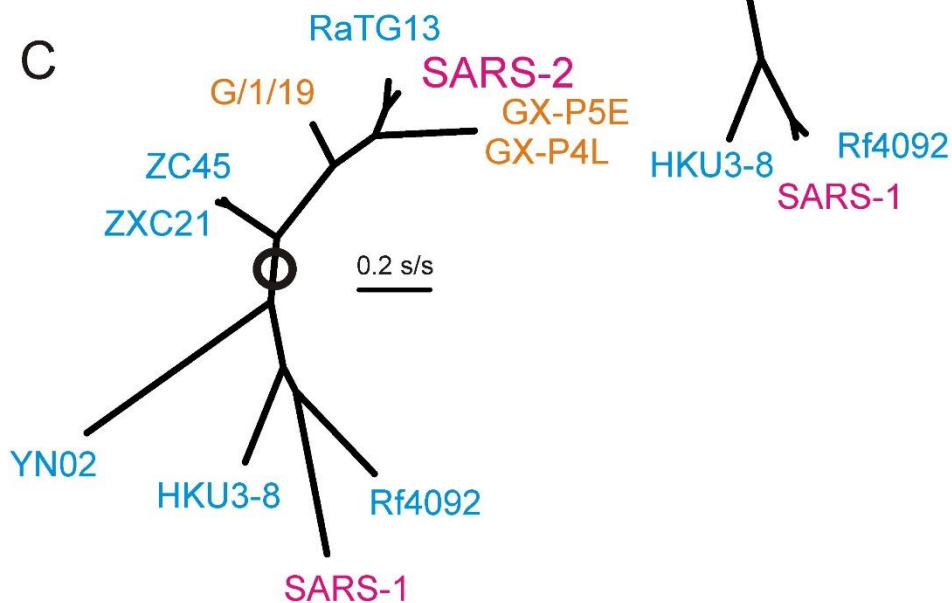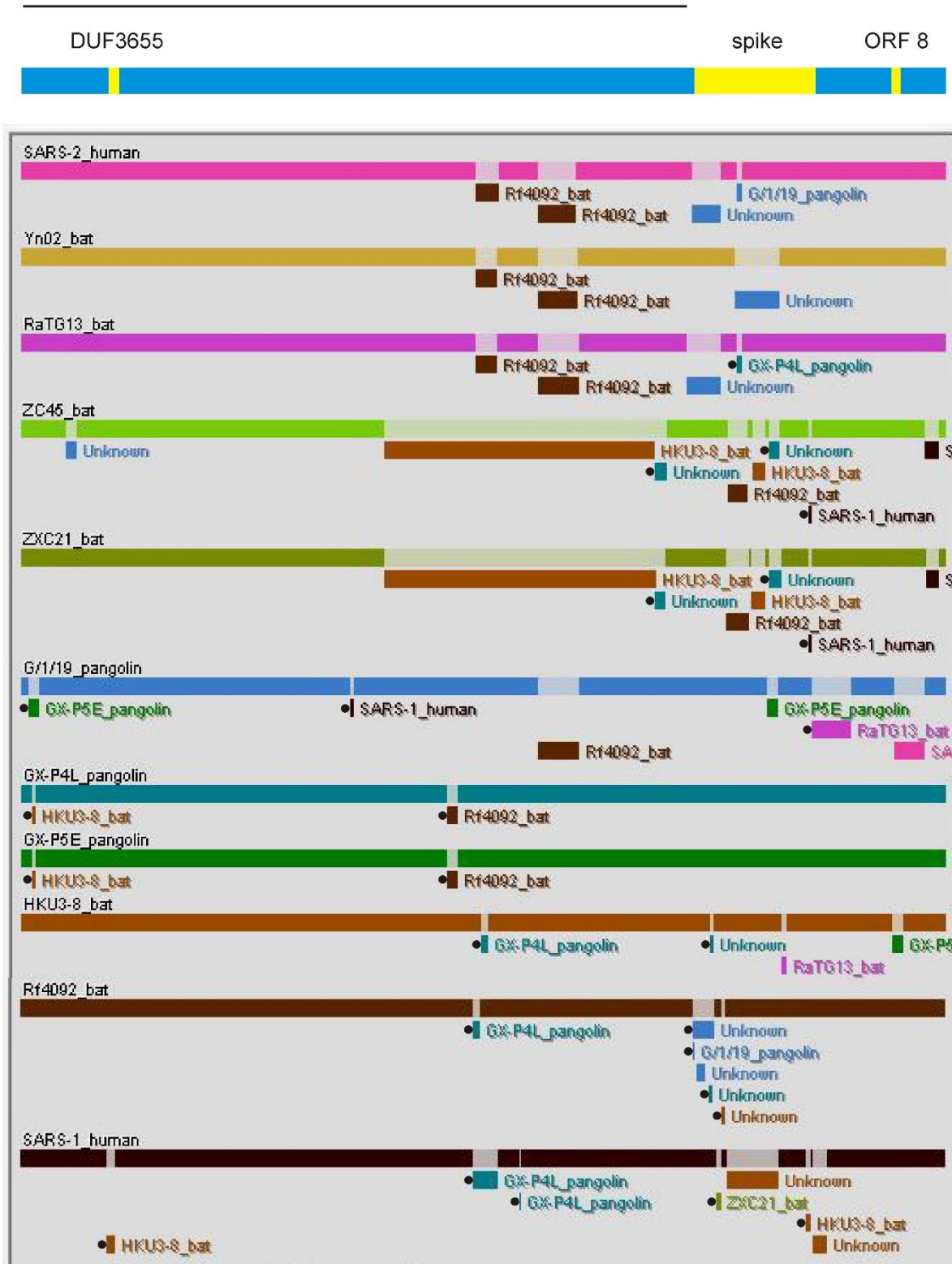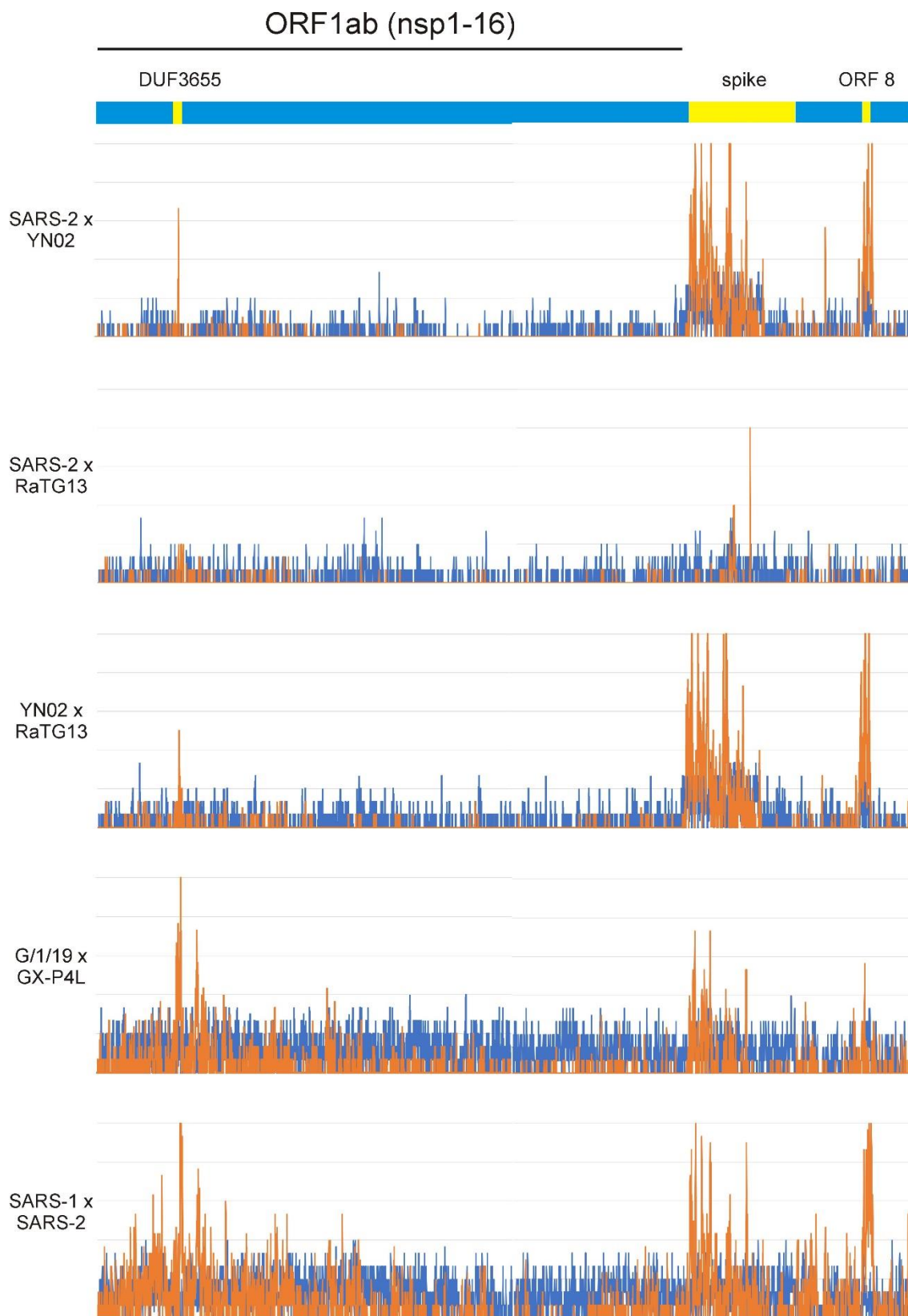
**Fig. 1**

**Fig. 2**

**Fig. 3**

## References

Abascal, F, Zardoya, R. and Telford, MJ., 2010. TranslatorX: multiple alignment of nucleotide sequences guided by amino acid translations. Nucleic Acids Res. 38(2):W7-W13.

Altschul, SF, et al. 1990. Basic local alignment search tool. J. Mol. Biol. 215:403-410.

Andersen KG, et al. 2020 The proximal origin of SARS-CoV-2. Nat Med 26 (4):450-452. doi:10.1038/s41591-020-0820-9

Angeletti, S. et al. 2020 The role of the nsp2 and nsp3 in its pathogenesis J Med Virol. 92:584–588. DOI: 10.1002/jmv.25719

Anthony SJ, et al Consortium P 2017 Global patterns in coronavirus diversity. Virus Evol 3 (1):vex012. doi:10.1093/ve/vex012

Boni, M. F., Posada, D., & Feldman, M. W. (2007). An exact nonparametric method for inferring mosaic structure in sequence triplets. Genetics, 176(2), 1035-1047.

Buchan DWA, Jones DT 2019 The PSIPRED Protein Analysis Workbench: 20 years on. Nucleic Acids Research. https://doi.org/10.1093/nar/gkz297

Candido et al 2020, Evolution and epidemic spread of SARS-CoV-2 in Brazil: https://www.medrxiv.org/content/10.1101/2020.06.11.20128249

Chan, YA. and Zhan, SH. 2020 Single source of pangolin CoVs with a near identical Spike RBD to SARS-CoV-22. bioRxiv preprint doi: https://doi.org/10.1101/2020.07.07.184374.

Chen L, Zhong L, 2020 Genomics functional analysis and drug screening of SARS-CoV-2, Genes & Diseases, https:// doi.org/10.1016/j.gendis.2020.04.002

Cozzetto D, et al. 2016 FFPred 3: feature-based function prediction for all Gene Ontology domains. Sci Rep. 2016 Aug 26;6:31865. doi: 10.1038/srep31865

Desikan, R. 2020 An In silico Algorithm for Identifying Amino Acids that Stabilize Oligomeric Membrane-Toxin Pores through Electrostatic Interactions bioRxiv preprint doi: https://doi.org/10.1101/716969.

Duchêne et al 2020, Temporal signal and the phylodynamic threshold of SARS-CoV-2. https://www.biorxiv.org/content/10.1101/2020.05.04.077735v1.full

Duchêne S, Holmes EC, Ho SY 2014 Analyses of evolutionary dynamics in viruses are hindered by a time-dependent bias in rate estimates. Proc Biol Sci 281, 1786. doi:10.1098/rspb.2014.0732

Ferron F, et al. 2018 Structural and molecular basis of mismatch correction and ribavirin excision from coronavirus RNA. Proc Natl Acad Sci U S A. 115(2):E162-E171. doi: 10.1073/pnas.1718806115. Epub 2017 Dec 26.

Fourment, M. and Gibbs, M.J., 2006. PATRISTIC: a program for calculating patristic distances and graphically comparing the components of genetic change. BMC evolutionary biology, 6(1), 1-5.

Gasteiger, E. et al 2005. Protein identification and analysis tools on the ExPASy server. In The proteomics protocols handbook (pp. 571-607). Humana press.

Gibbs MJ, et al 2007 The variable codons of H3 influenza A virus haemagglutinin genes. Arch Virol 152 (1):11-24. doi:10.1007/s00705-006-0834-8

Gibbs, MJ., Armstrong, JS., and Gibbs, AJ., 2000. Sister-Scanning: a Monte Carlo procedure for assessing signals in recombinant sequences. Bioinformatics 16:573-582.

Gorbalenya AE, et al Coronaviridae Study Group of the International Committee on Taxonomy of V (2020) The species Severe acute respiratory syndrome-related coronavirus: classifying 2019-nCoV and naming it SARS-CoV-2. Nature Microbiology 5 (4):536-544. doi:10.1038/s41564-020-0695-z

Gribble, J. et al 2020 The coronavirus proofreading exoribonuclease mediates extensive viral recombination bioRxiv preprint doi: https://doi.org/10.1101/2020.04.23.057786.

Guindon, S. and Gascuel, O., 2003. A simple, fast, and accurate algorithm to estimate large phylogenies by maximum likelihood. Systematic Biol. 52:696-704.

Hall, T. A., 1999. BioEdit: a user-friendly biological sequence alignment editor and analysis program for Windows 95/98/NT. Nucleic Acids Symp. Ser. 41:95-98.

Holmes, E. C., Worobey, M., & Rambaut, A. (1999). Phylogenetic evidence for recombination in dengue virus. Molecular biology and evolution, 16(3), 405-409.

Hsin, W. et al 2018 Nucleocapsid protein-dependent assembly of the RNA packaging signal of Middle East respiratory syndrome coronavirus. J Biomed Sci 25, 47 (2018). https://doi.org/10.1186/s12929-018-0449-x

Jeanmougin, F. et al 1998. Multiple sequence alignment with Clustal X. Trends Biochem. Sci. 23:403-405.

Katoh, K., and Standley, D. M., 2013. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. Mol. Biol. Evol. 30:772-780.

Latinne, A. et al 2020 Origin and cross-species transmission of bat coronaviruses in China. bioRxiv preprint doi: https://doi.org/10.1101/2020.05.31.116061.

Lemey, P. et al 2009. Identifying recombinants in human and primate immunodeficiency virus sequence alignments using quartet scanning. BMC Bioinformatics 10:126.

Li, X. et al 2020 Emergence of SARS-CoV-2 through Recombination and Strong Purifying Selection Short Title: Recombination and origin of SARS-CoV-2 bioRxiv preprint doi: https://doi.org/10.1101/2020.03.20.000885.t

Lin, X. and Chen, S. 2020 Major Concerns on the Identification of Bat Coronavirus Strain RaTG13 and Quality of Related Nature Paper.Preprints, 2020060044 (doi: 10.20944/preprints202006.0044.v1

Martin DP et al. 2015 RDP4: Detection and analysis of recombination patterns in virus genomes. Virus Evol 1 (1):vev003. doi:10.1093/ve/vev003

Martin, DP., and Rybicki, E., 2000. RDP: detection of recombination amongst aligned sequences. Bioinformatics 16:562-563.

Martin, DP et al., 2005. A modified bootscan algorithm for automated identification of recombinant sequences and recombination breakpoints. AIDS Res. Human Retroviruses 21:98-102.

Maynard-Smith, J.M., 1992. Analyzing the mosaic structure of genes. J. Mol. Evol. 34:126–129.

McBride R. van Zyl, M. Fielding, BC 2014 Review The Coronavirus Nucleocapsid Is a Multifunctional Protein. Viruses 6, 2991-3018; doi:10.3390/v6082991

McGuire, G., and Wright, F., 2000. TOPAL 2.0: Improved detection of mosaic sequences within multiple alignments. Bioinformatics 16:130–134.

Michalska, K. et al 2020 Crystal structures of SARS-CoV-2 ADP-ribose phosphatase (ADRP): from the apo form to ligand complexes. bioRxiv preprint doi: https://doi.org/10.1101/2020.05.14.096081.

Padidam, M., Sawyer, S., & Fauquet, C. M. (1999). Possible emergence of new geminiviruses by frequent recombination. Virology, 265(2), 218-225.

Posada, D. and Crandall, K. A., 2001. Evaluation of methods for detecting recombination from DNA sequences: computer simulations. PNAS 98:13757-13762.

Prates ET, et al., 2020 Functional Immune Deficiency Syndrome via Intestinal Infection in COVID-19. bioRxiv Confronting the COVID-19 Pandemic with Systems Biology bioRxiv preprint doi: https://doi.org/10.1101/2020.04.06.028712. t,

Pybus, OG et al 2020, Preliminary analysis of SARS-CoV-2 importation & establishment of UK transmission lineages: "E " (https://virological.org/t/preliminary-analysis-of-sars-cov-2-importation-establishment-of-uk-transmission-lineages/507/2).

Resende et al 2020, Genomic surveillance of SARS-CoV-2 reveals community transmission of a major lineage during the early pandemic phase in Brazil: " (https://virological.org/t/genomic-surveillance-of-sars-cov-2-reveals-community-transmission-of-a-major-lineage-during-the-early-pandemic-phase-in-brazil/514).

Rodríguez-Román E, Gibbs AJ. Ecology and Evolution of Betacoronaviruses. In: Rezaei N, ed. Coronavirus disease (COVID-19). Springer Nature 2020. In press.

Schmid K, Yang Z 2008 The trouble with sliding windows and the selective pressure in BRCA1. PLoS ONE 3(11): e3746. doi:10.1371/ journal.pone.0003746

Shimodaira, H., & Hasegawa, M. (1999). Multiple comparisons of log-likelihoods with applications to phylogenetic inference. Molecular biology and evolution, 16(8), 1114-1114.

Verheije, MH 2010 The Coronavirus Nucleocapsid Protein Is Dynamically Associated with the Replication-Transcription Complexes J. Virol. 84. 11575–11579

Wang, H., Pipes, L. Nielsen, R. (2020) Synonymous mutations and the molecular evolution of SARS-Cov-2 origins bioRxiv preprint doi: https://doi.org/10.1101/2020.04.20.052019

Wei, J. et al. 2020 Genome-wide CRISPR screen reveals host genes that regulate SARS-CoV-2 infection bioRxiv preprint doi: https://doi.org/10.1101/2020.06.16.155101.

Wolff, G. 2020 A molecular pore spans the double membrane of the coronavirus replication organelle bioRxiv preprint doi: https://doi.org/10.1101/2020.06.25.171686.

Wu, F., Zhao, S., Yu, B.et al.A new coronavirus associated with human respiratory disease in China.Nature 579, 265–269 (2020). https://doi.org/10.1038/s41586-020-2008

Xiao, K. et al. 2020 Isolation of SARS-CoV-2-related coronavirus from Malayan pangolins. Nature (2020) doi:10.1038/s41586-020-2313-x.

Ye, W. et al 2020 Zoonotic origins of human coronaviruses Int. J. Biol. Sci. 16(10): 1686-1697. doi: 10.7150/ijbs.45472

Zhou P et al. 2020 A pneumonia outbreak associated with a new coronavirus of probable bat origin. Nature 579 (7798):270-273. doi:10.1038/s41586-020-2012-7