Proposed title:
**Identification of evolutionary trajectories shared across human betacoronaviruses**

Authors listed according to affiliation:
1  Marina Escalera-Zamudio [1,2*]
2  Sergei L. Kosakovsky Pond [3]
3  Natalia Martínez de la Viña [1]
4  Bernardo Gutiérrez [1,2]
5  Rhys P. D. Inward [1]
6  Julien Thézé [4]
7  Lucy van Dorp [5]
8  Hugo G. Castelán-Sánchez [2,6]
9  Thomas A. Bowden [7]
10  Oliver G. Pybus [1,8]
11  Ruben J.G. Hulswit [7*]
12
13  Affiliations:
14      1.  Department of Biology, University of Oxford, Oxford, OX1 3PS, UK
15      2.  Consorcio Mexicano de Vigilancia Genómica (CoViGen-Mex)
16      3.  Institute for Genomics and Evolutionary Medicine, Temple University, Philadelphia, PA 19122, USA
17      4.  Université Clermont Auvergne, INRAE, VetAgro Sup, UMR EPIA, 63122, Saint-Genès-Champanelle, France
18      5.  UCL Genetics Institute, Department of Genetics, Evolution and Environment, University College London, Gower Street,
19          WC1E 6BT, London
20      6.  Programa de Investigadoras e Investigadores por México, Consejo Nacional de Ciencia y Tecnología, CP 03940, CDMX,
21          México
22      7.  Division of Structural Biology, Wellcome Centre for Human Genetics, University of Oxford, Oxford, OX3 7BN, UK
23      8.  Department of Pathobiology, Royal Veterinary College, NW1 0TU, London, United Kingdom
24
25  Email addresses (*corresponding authors):
26  marina.escalerazamudio@biology.ox.ac.uk *
27  spond@temple.edu
28  nataliamv@gmail.com
29  bernardo.gutierrez@biology.ox.ac.uk
30  rhys.inward@biology.ox.ac.uk
31  julien.theze@inrae.fr
32  lucy.dorp.12@ucl.ac.uk
33  hugo.castelan@conacyt.mx
34  thomas.bowden@strubi.ox.ac.uk
35  oliver.pybus@biology.ox.ac.uk
36  ruben.hulswit@strubi.ox.ac.uk *
37

38  **Preprint available at:** https://www.biorxiv.org/content/10.1101/2021.05.24.445313v2.full
39
40

41  **ABSTRACT (239)**

42  Comparing the evolution of distantly related viruses can provide insights into common adaptive
43  processes related to shared ecological niches. Phylogenetic approaches, coupled with other
44  molecular evolution tools, can help identify mutations informative on adaptation, whilst the
45  structural contextualization of these to functional sites of proteins may help gain insight into their
46  biological properties. Two zoonotic betacoronaviruses capable of sustained human-to-human
47  transmission have caused pandemics in recent times (SARS-CoV-1 and SARS-CoV-2), whilst a
48  third virus (MERS-CoV) is responsible for sporadic outbreaks linked to animal infections.

1

1 Moreover, two other betacoronaviruses have circulated endemically in humans for decades
2 (HKU1 and OC43). To search for evidence of adaptive convergence between established and
3 emerging betacoronaviruses capable of sustained human-to-human transmission (HKU1,
4 OC43, SARS-CoV-1 and SARS-CoV-2), we developed a methodological pipeline to classify
5 shared non-synonymous mutations as putatively denoting homoplasy (repeated mutations that
6 do not share direct common ancestry) or stepwise evolution (sequential mutations leading
7 towards a novel genotype). In parallel, we look for evidence of positive selection, and draw upon
8 protein structure data to identify potential biological implications. We find 30 candidate
9 mutations, from which four [codon sites 18121 (nsp14/residue 28), 21623 (spike/21), 21635
10 (spike/25) and 23948 (spike/796); SARS-CoV-2 genome numbering] further display evolution
11 under positive selection and proximity to functional protein regions. Our findings shed light on
12 potential mechanisms underlying betacoronavirus adaptation to the human host and pinpoint
13 common mutational pathways that may occur during establishment of human endemicity.
14

15 **KEYWORDS**
16 Molecular Evolution, Phylogenomics, Convergence, Adaptation, Betacoronaviruses
17

18 **SIGNIFICANCE STATEMENT (115)**
19 Identifying adaptive convergence is intimately linked to the possibility of predicting evolutionary
20 trajectories in viruses relevant to global health. In this light, we undertook a comparative
21 approach to find evidence of adaptive convergence across betacoronaviruses capable of a
22 sustained human-to-human transmission (HKU1, OC43, SARS-CoV-1 and SARS-CoV-2). Our
23 methodology involved the development of a pipeline used for identifying mutations putatively
24 denoting homoplasy and or stepwise evolution that are also evolving under positive selection,
25 and with potential biological implications drawn from protein structural data. Coupled with future
26 experimental data and ongoing genomic surveillance, our results raise the possibility of
27 predicting how the evolutionary trajectory for SARS-CoV-2 may develop as the virus establishes
28 itself as endemic to humans.
29

30 **MAIN TEXT (6183)**

31 **INTRODUCTION (598)**

32 Understanding the mutational processes that lead to adaptation in RNA viruses is crucial for

33 developing effective control strategies. Due to their high mutation rates and small genomes,

34 RNA viruses often display rapid evolution. However, the vast majority of mutations are either

35 purged through purifying selection or are selectively neutral (Loewe and Hill 2010). Only a small

36 proportion of these may contribute to adaptive evolution and be consequently fixed through

37 positive selection (Ohta 1973; Pond et al. 2012). For most viral genomes, the mutational

38 pathways leading to adaptation are further constrained by functional and evolutionary

39 limitations, such as epistasis, which refers to the adaptive dependence of a given mutation on

40 the genetic background in which it appears (Dolan et al. 2018). Therefore, viral evolutionary

41 trajectories are often limited and may exhibit recurrent mutational patterns indicative of adaptive

2

1 convergence, especially when applied to independent virus populations that share ecological

2 niches (Gutierrez, Escalera-Zamudio, and Pybus 2019).

3       The OC43 and HKU1 embecoviruses and the SARS-CoV and SARS-CoV-2

4 sarbecoviruses are four betacoronavirus species capable of sustained human-to-human

5 transmission. OC43 and HKU1 were introduced into the human population through independent

6 zoonotic events estimated to have occurred at least 50 years ago and are associated with mild

7 respiratory disease (Su et al. 2016). In contrast, SARS-CoV and SARS-CoV-2 were

8 independently introduced more recently, causing severe pandemic outbreaks (W. Li et al. 2005;

9 Vijaykrishna et al. 2007; Andersen et al. 2020; Boni et al. 2020; Banerjee et al. 2021). In 2002,

10 SARS-CoV spread to more than 20 countries, causing a short-lived outbreak characterized by

11 sustained human-to-human transmission (Cheng et al. 2007). Although its circulation was

12 eventually halted, the virus displayed evidence of adaptation to the human population (He et al.

13 2004). Almost two decades later, SARS-CoV-2 spread globally, resulting in the current

14 pandemic, despite a low rate of adaptive change recorded during the early stages of the

15 outbreak (van Dorp et al. 2020; MacLean et al. 2021). The continuous circulation of OC43 and

16 HKU1 within the human population at a global scale has been accompanied by ongoing host-

17 specific adaptation. This is now also evident for SARS-CoV-2, exemplified by the constant

18 emergence of novel virus lineages across time and space, with sub-lineages now reflecting

19 regional endemic patterns (O'Toole et al. 2021).

20       As SARS-CoV-2 becomes established in humans, it will continue to adapt to overcome

21 the selective pressures exerted by the collective immune response of the human population

22 (Kissler et al. 2020). We hypothesize that adaptive convergence may occur across distantly

23 related betacoronaviruses circulating within the same ecological niche, specifically the human

24 host. To test this, we undertook a comparative analysis to search for evidence of shared

25 mutational pathways between established human-endemic embecoviruses and emerging

26 sarbecoviruses, with a focus on emerging mutations observed in SARS-CoV-2. We developed a

27 methodological pipeline that allows for the identification of non-synonymous mutations

28 (rendering amino acid substitutions) likely associated with adaptive convergence across multiple

29 virus species. Firstly, we detected amino acid substitutions shared across virus taxa, displaying

30 putative evidence of homoplasy or stepwise evolution. Secondly, we assessed whether these

31 substitutions were positively selected, and contextualized their location to functional regions of

32 viral proteins. Following our pipeline, we initially detected 30 candidate amino acid substitutions

33 displaying evolutionary patterns denoting putative homoplasy and/or stepwise evolution. We

34 subsequently identified four of these (sites 18121 [nsp14/27], 21623 [spike/21], 21635

3

1 [spike/25], and 23948 [spike/796], in SARS-CoV-2 genome coordinates) as positively selected,

2 and proximal to functional surfaces in nsp14 (Ma et al. 2015) and the spike (S) protein. Our

3 results provide a molecular-level context for common evolutionary trajectories that

4 betacoronaviruses may undergo during their adaptation to the human host.

5

6 **RESULTS (2364)**

7 <u>Patterns of genetic variability observed in human-infecting betacoronaviruses</u>

8 We performed phylogenetic analyses of human-infecting betacoronaviruses using an alignment

9 of the Orf1ab and S viral genes (see Methods section 1 and 2). The tree shown in Figure 1

10 provides a comprehensive picture of the evolutionary relationships among the four

11 betacoronavirus species studied here, consistent with previously published phylogenies of the

12 genus (Woo et al. 2006; Woo et al. 2010; Oong et al. 2017; Zhu et al. 2018; Bedford 2021). Our

13 analysis confirms four well-supported clades formed by virus sequences belonging to the

14 *Embecovirus* (HKU1, OC43, and related viruses) and *Sarbecovirus* (SARS-CoV, SARS-CoV-2,

15 and related viruses) subgenera (ICTV et al. 2017). To further validate divergence patterns at a

16 deeper node level, we compared individual clades (sub-trees within our trees) to species-

17 specific phylogenies. We were also able to verify the divergence patterns described for the

18 distinct HKU1 (A-C) and OC43 (A-H) genotypes (Woo et al. 2006; Oong et al. 2017)

19 (Supplementary Data 1). Therefore, our phylogenetic reconstructions validate the evolutionary

20 relationships among these four distantly related betacoronaviruses.

21 We then analysed the proportion of codon sites (from the total number of polymorphic

22 sites identified), corresponding to non-synonymous mutations shared between different

23 embeco- and sarbecovirus species (*i.e.,* those present in any of the sarbecovirus clades, and

24 also in HKU1 and/or OC43). Derived from the Orf1ab+S alignment (comprising a total of 8962

25 sites), we identify approximately 2% (205 sites) as shared. Within Orf1a region (4774 sites),

26 2.7% of these (129 sites) were identified as shared. Within Orf1b region (2623 sites), only 0.9%

27 (25 sites) were further identified as shared. The Orf S region (1457 sites) displayed the highest

28 proportion of shared mutations (3.2%, 48 sites). When analysing genetic variation patterns

29 within single virus species, we observed a high degree of sequence conservation (>91%

30 identity) across the Orf S of all virus species. Conserved sites were predominantly located in the

31 membrane-proximal S2 domain, while variable sites were mostly found within the membrane-

32 distal S1 subunit (Figure 2). The predominance of variable sites within S1 compared to S2 was

33 most evident for embecoviruses, and less so for sarbecoviruses, suggesting for a differential

4

1   adaptation stage relative to the human host environment, evidenced by a lower degree of

2   genetic divergence observed in Orf S in the sarbecoviruses.

3         We further analysed the genetic variation across virus species, focusing on the Orf S

4   region. As previously noted for other coronaviruses (Hulswit et al. 2016), we found that Orf S

5   exhibited a higher proportion of variable sites relative to conserved (for definitions, see Methods

6   section 3). Specifically, only 16% of homologous sites within the Orf S alignment were

7   conserved, while the remaining 84% were variable (Supplementary Data 2). The S2 subunit of

8   Orf S contained the highest proportion of conserved sites, presumably due to shared functional

9   constraints of the viral membrane fusion machinery across coronavirus species (Li 2016).

10  Conversely, the S1 subunit displayed a higher number of variable sites, particularly within the

11  S1$^A$ domain (also known as the N-terminal domain, or NTD). We found that the S1$^B$ domain did

12  not display any conserved sites across virus species, likely due to differences in receptor

13  engagement between embeco- and sarbecoviruses. Specifically, embecoviruses use the S1$^A$

14  domain to interact with sialoglycan-based receptors, while sarbecoviruses use their S1$^B$ domain

15  to bind to angiotensin-converting enzyme 2 (ACE2) (Hulswit et al. 2019; Lan et al. 2020).

16  Finally, we identified that the conserved R residue at site 685 corresponding to the S1/S2

17  cleavage site (numbering according to the SARS-CoV-2 protein, codon sites 23615-23617) is

18  shared across and within virus species (Supplementary Data 2), reflecting a conserved

19  proteolytic maturation mechanism of the spike protein (Millet and Whittaker 2015).

20

21  <u>Sites displaying evidence of homoplasy and/or stepwise evolution</u>

22  Although not all non-synonymous mutations putatively displaying homoplasy and/or stepwise

23  evolution may arise from positive section, such mutational patterns are most likely to result from

24  adaption (Escalera-Zamudio et al. 2020; Stern et al. 2017; Gutierrez et al. 2019). Thus, amongst

25  the non-synonymous mutations identified as shared across virus species, we further searched

26  for those displaying putative evidence for homoplasy and/or stepwise evolution (Supplementary

27  Text 1) using our pipeline (Methods section 3). After visual validation, we confirmed that 30 sites

28  (representing 0.3% within the Orf1ab+S alignment) display evolutionary patterns indicative of

29  homoplasy and/or stepwise evolution (see Supplementary Text 2, Supplementary Figure 2 and

30  3). Two of these were found within Orf1a, nine within Orf1b, and 19 within Orf S (Table 1). The

31  evolutionary trajectories for different amino acid states observed for three illustrative sites

32  (18121, 21623 and 23948, further displaying evidence of evolution under positive selection and

33  of being proximal to regions of established protein function [see the following results sections])

5

are highlighted below (Figure 3). The amino acid evolution patterns observed for all other sites are available in Supplementary Data 3.

Derived from the global, expanded and the re-sampled SARS-CoV-2 trees (Methods section 1, 3 and 6, Supplementary Text 3), our results show that site 18121 (codon 18121-18123 in Orf1b, corresponding to amino acid state 'S' in nsp14 in SARS-CoV-2 numbering) is homoplasic between HKU1 genotype B and the sarbecoviruses (Table 1, Figure 3, Supplementary Data 3). Comparably, site 21623 (codon 21623-21626 in Orf S, corresponding to amino acid state 'R' in S) was identified as homoplasic between SARS-CoV-2 and OC43 genotypes D, F, G and H. This site also displayed evidence for stepwise evolution within a single virus clade (OC43), exemplified by the sequential amino acid replacement pattern of V→ I→ K→ R (Figure 3).

For site 23948 (codon 23948-23950 in Orf S, corresponding to residue 796 in S), initial observations based on the global tree revealed that amino acid state 'D' was present in all virus species, except for OC43 (displaying amino acid state 'N'). However, when replicating our analyses (expanded tree), the distribution of amino acid state 'D' was now found present in some embecoviruses (including OC43 but excluding HKU1) and most sarbecoviruses. These discrepancies are likely due to alignment uncertainty across genome regions of highly divergent virus taxa. Nonetheless, based on consensus protein sequences and structural comparison, the structural contextualization of amino acid 796 and adjacent sites confirmed the presence of 'D' in SARS-CoV-1, SARS-CoV-2 and HKU1, and 'N' in OC43, (Supplementary Figure 6). Thus, amino acid state 'D' at site 23948 shows evidence of homoplasy between the SARS-CoV-1, SARS-CoV-2 and HKU1.

For this same site (23948), an additional amino acid change from 'D' to 'Y' was identified as homoplasic between some SARS-CoV-1 and SARS-CoV-2 sequences (data derived from the global, expanded and the re-sampled SARS-CoV-2 trees) (Table 1). For SARS-CoV-2, amino acid state 'Y' emerged and was lost repeatedly during the early stage of the pandemic (represented by independent minor clusters that quickly became extinct). However, following emergence and global spread of the B.1.1.529 virus lineage (Omicron variant of concern [VOC], and descending sub-lineages), amino acid state 'Y' replaced amino acid state 'D', displaying a predominant trend associated with the dominance of the B.1.617.2 lineage (Delta VOC, and descending sub-lineages) (Table 1, Figure 3, Supplementary Data 3) (also confirmed by independently sampled SARS-CoV-2 data available up to December 2022: https://nextstrain.org/groups/neherlab/ncov/global?c=gt-S_796).

6

1  Quantifying the effects of positive selection

2  The dN/dS estimates we obtained across complete virus genomes and upon specific coding

3  regions (see Methods section 4) indicates that positive selection is acting upon the Orf1ab and

4  Orf S of SARS-CoV-2, compared to other viruses studied here. Specifically, the effect of

5  episodic diversifying selection was detected upon 5/14 non-recombinant fragments (three in

6  Orf1b and two in Orf S, for details see https://observablehq.com/@spond/beta-cov-analysis).

7  Using the Contrast-FEL method to detect the effect of a differential selection across branches

8  separating lineages (see Methods section 4), we found 36 sites (0.4%) evolving under

9  differential selective pressure across distinct virus clades. Furthermore, we found 0.7% of all

10  sites (67 codons within the Orf1ab+S alignment) to be evolving under episodic diversifying

11  positive selection (scored under MEME with a p≤0.05 as *positively-selected sites*, or PSS)

12  (Supplementary Table 1). In contrast, we found 5% of all sites (461 codons within the Orf1ab+S

13  alignment) to be evolving under pervasive negative selection (scored under FEL with a p≤0.05

14  as *negatively-selected sites,* or NSS). We subsequently mapped the identified PSS and NSSs

15  onto the SARS-CoV-2 S protein structure (Methods section 5). We observe that out of a total of

16  22 PSSs, 18 locate within the S1 subunit (11 in S1$^A$, 5 in S1$^B$, 1 in S1$^C$ and 1 in S1$^D$ domains),

17  whilst the remaining four mapped onto the S2 subunit. Conversely, out of a total of 82 of NSSs,

18  46 locate within S1 (18 in S1$^A$, 21 in S1$^B$, 3 in S1$^C$ and 4 in S1$^D$), whilst the remaining 36

19  mapped onto S2 (Supplementary Figure 7).

20  From the 30 non-synonymous mutations we identify as displaying evolutionary patterns

21  putatively denoting homoplasy and/or stepwise evolution (Table 1), sites 19048, 21623, 21635,

22  22124 and 23048 were further scored as PSS (under different methods). Sites 21623 and

23  21635 were inferred as PSSs along ancestral branches leading to the HKU1, OC43 and SARS-

24  CoV-2 clades. Sites 19048 and 22124 were inferred as PSSs along the OC43 ancestral branch,

25  whilst 23048 was inferred as a PSS along the HKU1 ancestral branch (Table 1, Supplementary

26  Table 1). Further analysis under the branch and site model in the MEME method (Methods

27  section 4) revealed site 18121 to be evolving under positive selection for the HKU1

28  clade/branch (relative to the sarbecoviruses), in agreement with our observations made on

29  putative homoplasy detected for this site between HKU1 genotype B, SARS-CoV-1 and SARS-

30  CoV-2 (Table 1, Figure 3). Similarly, site 23948 was also inferred to be evolving under positive

31  selective for the SARS-CoV-1 branch, relative to other virus clades (Supplementary Table 1).

32  For validation, we compared our results with selection analysis available for

33  independently sampled SARS-CoV-2 genome data available as of December 2022

34  (https://observablehq.com/@spond/evolutionary-annotation-of-sars-cov-2-covid-19-genomes-

7

1  [enab](#)) (Kosakovsky Pond). Of the 30 mutations we identify, 16 of these are currently scored as

2  PSS or NSSs, with 13 of these mapping directly onto potential T cell epitopes derived from HLA

3  class I and HLA-DR binding peptides in SARS-CoV-2 (Nelde et al. 2021; Campbell et al. 2020)

4  (Table 1). Additionality, up to December 2022, sites 7478, 21614, 23948, 24620 and 25166

5  were detected as evolving under positive selection, whilst sites 21635, 24863, and 25037 were

6  detected as evolving under negative selection.

7

8  <u>Contextualization of mutations using protein structural and functional information</u>

9  We then mapped the 30 mutations identified onto corresponding protein structures. Below, we

10  focus on four exemplary sites (18121, 21623, 21635, and 23948) that meet the three criteria of:

11  displaying evidence of homoplasy and/or stepwise evolution, showing evidence of evolution

12  under positive selection, and being proximal to regions of established protein function. A

13  description for the other 26 identified mutations is available in the Supplementary Text 4 and

14  Supplementary Table 2.

15

16  **Site 18121 in Orf1ab**

17  Site 18121 is located within the Orf1ab gene and corresponds to an 'S' to 'A' mutation at residue

18  28 within the exonuclease domain of the nsp14 protein (numbering according to the SARS-CoV

19  protein) (Supplementary Table 2, Figure 4). Nsp14 is involved in the 5′-capping of viral mRNA

20  and is essential for viral mRNA transcription (Ma et al. 2015). The 'S' to 'A' mutation within this

21  region is expected to result in the loss of an intra-protein hydrogen bond and potentially

22  modulates    the    protein-protein    interaction    (Figure    4)    (assessed    under    PISAebi;

23  [http://www.ebi.ac.uk/pdbe/prot_int/pistart.html](http://www.ebi.ac.uk/pdbe/prot_int/pistart.html)) (Krissinel and Henrick 2007).

24

25  **Sites 21623 and 21635 in S1**

26  The S1 subunit mediates attachment of the virus to the host cell (Li 2016). Human-infecting

27  embecoviruses bind to glycan-based cell receptors via two hydrophobic pockets within the S1[A]

28  region of the protein protein (Hulswit et al. 2019; Tortorici et al. 2019), while the receptor-binding

29  site for human-infecting sarbecoviruses is located within the S1[B] domain of the protein (Li et al.

30  2005; Lan et al. 2020; Shang et al. 2020). Both SARS-CoV and SARS-CoV-2 recognize the

31  ACE2 molecule to enter the host cell, despite limited conservation amongst contact residues

32  within the RBD of these virus species (Li et al. 2005; Lan et al. 2020). Site 21623 displays

33  several non-synonymous mutations ('R', 'V', 'K' and 'I') mapping to residue 29 within the S1[A]

34  domain of the S1 subunit. Site 21635 also shows multiple non-synonymous mutations ('P', 'V',

8

1   'S', 'L' and 'H') mapping to residue 33 in S1[A]. For the OC43 S protein, this corresponds to a loop

2   neighbouring the hydrophobic pockets in S1[A] instrumental for receptor recognition (Figure 5),

3   and changes within this region may potentially modulate receptor affinity (Hulswit et al. 2019).

4   The mutational patterns observed at these sites putatively denote homoplasy/stepwise evolution

5   and evidence of positive selection (Table 1), and are therefore congruent with antigenic drift

6   shaping the evolution of human-endemic coronaviruses (Kistler and Bedford 2021). In SARS-

7   CoV-2, mutations in both these sites (residue 29 and 33) have been observed for two VOCs

8   (B.1.351 and P.1, 'Beta' and 'Gamma') (Faria et al. 2021; Tegally et al. 2021). Even though

9   sarbecoviruses engage the ACE2 receptor via domain S1[B], these residues locate to the 'NTD

10  supersite', serving as epitope for multiple of neutralizing antibodies (Kemp et al. 2021).

11

## Site 23948 in S2

13  The S2 subunit of the betacoronavirus S protein contains the fusion machinery, responsible for

14  merging the viral envelope with the host cell membrane to facilitate delivery of the viral genome

15  into the target cell. This process is driven by the fusion peptide, which anchors the virus to the

16  host membrane, and requires cleavage of the S protein by host cell proteases at the S1-S2

17  junction (consensus RRAR|S in SARS-CoV-2) and at the S2' cleavage site (R|S, located

18  immediately upstream of the fusion peptide in the S2 subunit) (Li 2016, Millet and Whittaker

19  2015). Site 23948 displays a non-conservative amino acid replacement 'D' to 'Y' (identified as

20  homoplasic between some SARS-CoV and SARS-CoV-2 sequences) at residue 796 of the S2

21  subunit, located immediately upstream of the S2' cleavage site (Table 1, Supplementary Table

22  2). This residue locates within a loop crucial for the release of the fusion peptide, exhibiting

23  some variability across betacoronavirus species (Supplementary Figure 6). Our observations

24  suggest that the apparent relaxed local constraints at this site may facilitate cleavage activation

25  by securing loop accessibility. Perhaps consistently, the corresponding protein region in the

26  HKU1 structure remains unresolved (Kirchdoerfer et al. 2016).

27

## DISCUSSION (1278)

29  In this study, we searched for signatures of adaptive convergence across distantly related

30  human-infecting betacoronaviruses, represented by shared non-synonymous mutations that

31  putatively denote homoplasy and/or stepwise evolution, further ranked according to their

32  selective relevance, and to their proximity to protein regions of known function. The majority of

33  the mutations we observe locate to the receptor binding region of the S protein (*i.e.*, S1 subunit),

34  whilst a smaller proportion of these were found within non-structural proteins encoded by Orf1ab

9

1 (site 18121 in the exonuclease domain of nsp14, and site 20344 in the endonuclease domain of

2 nsp15). Our *in-silico* analyses revealed four genomic sites (18121, 21623, 21635 and 23948)

3 that display cumulative evidence of: *i*) a mutational pattern putatively denoting homoplasy

4 and/or stepwise evolution, *ii*) evolution under positive selection, and *iii*) being structurally

5 proximal to regions of known protein function. Below, we discuss our findings in light of three

6 key evolutionary processes: antigenic drift, epistasis and adaptive convergence.

7       The host humoral immune response is an evolutionary force driving viral antigenic drift.

8 In the case of betacoronaviruses, this is reflected by cumulative mutations in the S protein

9 (particularly within the S1 subunit) that may allow frequent reinfections of the host population

10 (Kistler & Bedford 2021; Yewdell 2021; Forni et al. 2021). In agreement with this observation,

11 the emergence of some SARS-CoV-2 lineages (particularly VOC) has been associated with

12 high levels of infection in pre-exposed human populations across different geographic regions

13 (as an example on P.1, see Faria et al. 2021). Our results evidence antigenic drift upon the S1

14 subunit of distinct betacoronaviruses as a major component of the adaptation process to the

15 human host environment, further evidenced by Orf S also being the least conserved genome

16 region across distinct virus species (Li 2016). On the other hand, mutations found within Orf1ab

17 could have a potential impact on viral fitness related to an enhanced replication efficacy in the

18 human host (Menachery et al. 2017). As the evolution of Orf1ab is also driven by immune

19 responses such as cytokine signalling cascades and antigen presentation (Wang et al. 2015;

20 Taefehshokr et al. 2020; Hackbart et al. 2020; Yuen et al. 2020), these mutations may also be

21 the result of concerted selective pressure(s), following that single mutational changes can have

22 pleiotropic effects on distinct viral phenotypes and fitness components (de Wilde et al. 2018).

23       Identifying adaptive convergence raises the possibility of predicting mutational pathways

24 in viruses important to global health (Gutierrez et al. 2019). When applied to SARS-CoV-2, our

25 results reveal that some of the mutations we had initially identified as potentially relevant back in

26 May 2021 (see (Escalera-Zamudio et al. 2021) had already been observed in other

27 betacoronaviruses that circulate endemically in humans (Table 1), and some now display

28 dominant trends in SARS-CoV-2 (as analysed up to December 2022). For example, amino acid

29 state 'R' at residue 21 of the S protein (sites 21623)

30 (https://nextstrain.org/groups/neherlab/ncov/global?c=gt-S_21) and 'P' at residue 25 (site

31 21635) (https://nextstrain.org/groups/neherlab/ncov/global?c=gt-S_25) have dominated across

32 time. Moreover, mutation 'D' to 'Y' observed at residue 796 of the S protein (site 23948) has

33 proven to be a successful mutational pathway, evidenced the replacement of amino acid state

34 'D' (previously observed for the B.1.617.2 lineage, Delta VOC and descending sub-lineages) by

10

1    'Y' (now observed for the B.1.1.529 lineage, Omicron VOC and descending sub-lineages)
2    (https://nextstrain.org/groups/neherlab/ncov/global?c=gt-S_796). Of interest, mutations at
3    residue 796 of the S protein have been linked to the emergence viral variants that display
4    reduced susceptibility to neutralizing antibodies (Kemp et al. 2021).

5    Epistasis is thought to have played a central role in the emergence of human-infecting
6    betacoronaviruses (Holmes & Rambaut 2004). However, inferring epistasis across diverging
7    viruses is difficult given the functional differences between homologous genes and proteins.
8    Through our methodological approach we cannot measure epistasis *per se*, but we can aim to
9    identify adaptive convergence and subsequently discuss its possible effects. Thus, our results
10   indirectly provide support for epistasis, in the sense that if the same amino acid changes are
11   observed in different virus species, then associated epistatic interactions are expected to be
12   shared. This is of particular importance when considering the potential role of epistasis in
13   antigenic drift, where the combined effect of independent mutations could contribute to antigenic
14   escape (Rochman et al 2022). In the context of our findings, sites 21623 and 21635 are
15   presumed to be involved in the antigenic drift of embecoviruses. As these residues are in close
16   proximity to each other (displaying a linked evolution), these could thus reflect epistatic
17   interactions. Nevertheless, within the SARS-CoV-2 S1$^B$-ACE2 interface, epistasis seems to play
18   a limited role, as the effect of multiple mutations seems to be additive rather than epistatic
19   (Rochman et al. 2022; Zahradník et al. 2021; Starr et al. 2022).

20   The mutational spectrum of SARS-CoV-2 is known to be impacted by the human host
21   apolipoprotein B mRNA-editing enzyme (APOBEC) family (Di Giorgio et al. 2020). The activity
22   of APOBEC induces C → U/T mutations in the viral genome through a cytidine deaminase
23   activity, likely resulting in a high degree of apparent homoplasy reflected in emerging mutations
24   across distinct virus sub-populations (De Maio, et al. 2020; Worobey, et al. 2020; Wang, et al.
25   2021). Relative to more commonly used strategies for identifying homoplasy within single virus
26   species, our methodology poses an alternative approach that aims to identify homoplasy *across*
27   *and within* virus taxa, represented by shared mutations most likely fixed under an evolutionary
28   scenario driven by selection (see Supplementary Text 5). Given that candidate mutations are
29   observed over longer evolutionary times, this approach represents a useful tool to decrease the
30   likelihood of erroneously scoring mutations as homoplasic (such as those resulting from
31   mutational biases inherent to the SARS-CoV-2 genome evolution).

32   However, identifying adaptive convergence faces several important limitations. First, the
33   methodology we use is conservative, as it is based on strict homology. In this context, we only
34   consider sites robustly identifiable as homologous that can be traced back to ancestral nodes

11

1 with confidence (consequently excluding highly divergent genes). Therefore, our approach may
2 result in an underestimation of sites that may putatively denote adaptive convergence across
3 highly divergent viruses. Moreover, a limited virus genome sampling across time and space (in
4 particular for HKU1 and SARS-CoV-1), coupled with a relatively low genetic diversity observed
5 for SARS-CoV-2 (Rausch et al. 2020), further restricts the potential to identify shared mutations
6 across virus species (van Dorp et al. 2020). In addition, there is some uncertainty associated
7 with the mutations identified, as (though unlikely given cumulative evidence derived from
8 different methodological approaches) it is not possible to rule out that some of these may still
9 derive from biological processes other than adaptation (such as founder effects, mutational
10 hitch-hiking, linkage, and toggling at hypervariable sites) (Kosakovsky Pond et al. 2012; Delport
11 et al. 2008; De Maio et al. 2021; Wang et al. 2020; Simmonds 2020). Finally, whilst our analysis
12 provides insights into coronavirus evolution in humans, our approach renders us unable to
13 identify mutations that may result from host switching events. This is due to analyses on nodes
14 representing ancient host switching events (Corman et al. 2018) being constrained by long
15 divergence times, differences in mutation rates across virus taxa in different animal hosts,
16 mutational saturation, and by a considerable under-sampling of betacoronaviruses circulating in
17 non-human hosts (Holmes & Rambaut 2004; De Maio et al. 2021).

18 In this sense, additional/future experimental data could help reveal the impact of
19 mutations on viral fitness. However, performing such studies may be difficult, as these concern
20 potential gain-of-function experiments. Alternatively, enhanced genomic surveillance of
21 betacoronaviruses infecting the human population and of those ciruclating in other animal host
22 may confirm whether the mutational pathways we identify here represent evolutionary
23 trajectories on which betacoronaviruses converge in their adaptation process to the human host.
24

25 **MATERIAL AND METHODS (1933)**

26 **1. Initial data collation**

27 When this manuscript was first deposited as a preprint (May 2021) (Escalera-Zamudio et al.
28 2021), complete genomes for all HKU1, OC43 and SARS-CoV-1 viruses sampled across
29 different geographical regions and time were downloaded from the Virus Pathogen Resource
30 (ViPR-NCBI 2021) (Supplementary Data 4). Sequences were removed if meeting any of the
31 following criteria: (i) being >1000nt shorter than full genome length, (ii) being identical to any
32 other sequence, or (iii) if showing >10% of site were ambiguities (including N or X). A total of 53
33 HKU1, 136 OC43 and 40 SARS-CoV-1 sequences were initially retained for analyses. For
34 SARS-CoV-2, to better reflect an early zoonotic process into the human population (MacLean et

12

1  al. 2021), we originally aimed to limit the genetic diversity of the sampled virus population to the

2  first wave of infection recorded during the pandemic. For this, ~23000 full genomes sampled

3  worldwide before May 2021 available in the GISAID platform (GISAID 2021) were downloaded

4  and aligned as part of an initial public dataset provided by the COG-UK consortium (COG-UK

5  2021) (Supplementary Data 4). To make local analyses computationally feasible, the original

6  SARS-CoV-2 dataset was randomly subsampled to ~5% of its original size, keeping the earliest

7  genomes, and further reducing the dataset under the quality criteria stated above. In total, 1120

8  SARS-CoV-2 sequences were retained. For all virus species considered, we focused only on

9  genomes derived from human cases, in order to reflect host-specific adaptation processes.

10

## 11     2.  Phylogenetic analyses

12  Only the main viral ORFs (Orf1ab and S) were used for further phylogenetic analyses, as these

13  are homologous amongst the four viral species studied, and encode proteins essential to certain

14  stages of the virus life cycle (*i.e.,* replication and entry). For each virus species, individual ORFs

15  (codons) were extracted and aligned as translated amino acid sequences using MAFFT v7.471

16  (to be then reverted to codons again) (Katoh & Standley 2013). Individual alignments were

17  concatenated to further generate species-specific concatenated Orf1ab+S alignments. The

18  concatenated alignments were then combined to generate a global alignment comprising all

19  virus species, that was re-aligned again at an amino acid level using a profile-to-profile

20  approach following taxonomic relatedness (Wang & Dunbrack Jr 2004). The final alignment was

21  reverted to codon sequences as input for all further analyses. The global alignment comprised

22  in total 1314 sequences and 26883 sites.

23      Maximum Likelihood phylogenies were estimated for the individual and global codon

24  alignments using RAxML v8 (Stamatakis 2015), under a general time reversible nucleotide

25  substitution model and a gamma-distributed among-site rate variation (GTR+G). Branch support

26  was assessed using 100 bootstrap replicates. All trees were midpoint-rooted, whilst general

27  phylogenetic patterns observed amongst these distantly related virus species were validated by

28  comparing to previously published phylogenies (Woo et al. 2010; Zhu et al. 2018; Lau et al.

29  2011; Oong et al. 2017; Woo et al. 2006; Bedford 2021). Recombination is known to be

30  common amongst betacoronaviruses (Oong et al. 2017; Woo et al. 2006; Su et al. 2016),

31  including SARS-CoV-2 (Gutierrez et al. 2022; Turakhia et al. 2022). However, recombinant

32  sequences were not removed at this step, as it was important to detect potentially recombinant

33  isolates that could display relevant mutations. Putative recombinant sequences were eventually

34  removed for subsequent analyses (when identified, see Methods sections 6 and 7).

13

## 3. Identifying homoplasy and/or stepwise evolution

Following the pipeline described by Escalera & Golden (Escalera-Zamudio et al. 2020), variable sites across different virus taxa were identified within the global alignment as those displaying non-synonymous mutations (rendering amino acid changes) occurring in at least ≥1% of the sampled sequences. Variable sites were extracted by masking columns across the alignment showing identical sites and at least 50% gaps, followed by the 'Find Variations/SNPs' function used to compare each site to consensus sequences generated under a 95% threshold with Geneious Prime v2020.0.4 (Kearse et al. 2012). A total of 6681 variable sites were identified and used to infer ancestral amino acid state reconstructions onto the nodes/internal branches of the global tree (see Methods section 2 above). This was done using TreeTime (Sagulenko et al. 2018) under a ML approach (RAS-ML) using a time-reversible model (GTR) for state transitions. The genetic variability observed within leaves/tips of the tree was deliberately excluded, in order to only analyse changes occurring within nodes or internal branches. In parallel, conserved sites were identified as those present in ≥ 99% of the sampled virus sequences. Conserved sites were extracted by reversing the 'variable site masking', to obtain only identical sites identified across the global alignment (Supplementary Data 2).

The resulting 6681 'Ancestral Reconstruction Trees' (named here ARTs) were then classified under a computational algorithm developed to sort mutational patterns based on whether or not they support homoplasy and/or stepwise evolution. Briefly, homoplasy can occur within nodes of single clade or across clades, in which the same amino acid change must be present in at least one internal node of any given clade, and in another internal node of the same/another clade. Clades with the same amino acid states must not share direct common ancestry. Conversely, stepwise evolution is represented as sequential mutations occurring at the same sites within a single clade. Any given site scored under putative 'stepwise evolution' must display changes between at least two different states (A→B), but without any immediate reversion (B→A). A description of the definitions used here for homoplasy and/or stepwise evolution are available as Supplementary Text 1 and Supplementary Figure 1. A description of all basic steps used in our algorithm, including a schematic representation, is available in the Supplementary Text 2, Supplementary Figure 2 and 3. Associated code is publicly available at https://github.com/nataliamv/SARS-CoV-2-ARTs-Classification.

14

### 4. Estimating dN/dS

Derived from the global alignment and tree, we estimated dN/dS ($\omega$, the ratio between the non-synonymous substitution rate per non-synonymous site and the synonymous substitution rate per synonymous site) using the following site, branch and branch-site models: Mixed Effects Model of Evolution (MEME), Fixed Effects Likelihood (FEL), and the fixed effects site-level model (Contrast-FEL) (Kosakovsky Pond & Frost 2005; Kosakovsky Pond et al. 2021; Murrell et al. 2012). For this, the alignment was partitioned into 14 putatively non-recombinant regions using the Genetic Algorithm for Recombination Detection (GARD) (Kosakovsky Pond, Posada, et al. 2006), with all subsequent analyses conducted on the partitioned data. As dN/dS models use the GTR component for the nucleotide evolutionary rate, biased mutation rates are handled. Further, to mitigate the inflation in dN/dS estimates that results from unresolved and/or maladaptive evolution, testing for selection was again restricted to internal nodes/branches of the phylogeny (Kosakovsky Pond, Frost, et al. 2006). Genome-wide comparison of dN/dS estimates across viral genome regions was performed using the Branch-Site Unrestricted Statistical Test for Episodic Diversification method (BUSTED) (Murrell et al. 2015). Finally, the impact of changing biochemical properties at selected sites was further assessed under the Property Informed Models of Evolution method (PRIME) (HyPhy 2013). Our results were further compared to the selection analysis available for independently sampled SARS-CoV-2 genome data available as of December 2022 (https://observablehq.com/@spond/evolutionary-annotation-of-sars-cov-2-covid-19-genomes-enab) (Kosakovsky Pond).

### 5. Mapping mutations onto betacoronavirus protein structures

To locate the non-synonymous mutations identified on viral protein regions of known function, corresponding residues were mapped to available structural data using PyMOL v 2.4.0 (https://pymol.org/2/) (Supplementary Table 2, see Data Availability section). Mutations were analysed in the context of their relative proximity to previously reported functional regions, and to each other. N-linked glycosylation sites in S protein sequences were identified by searching for the N-[not P]-[S or T] consensus sequence (Watanabe et al. 2019). None of the mutations identified in this study resulted in generation or deletion of N-linked glycosylation sequons. In parallel, conserved and variable sites identified (including the 30 mutations evidencing homoplasy and/or stepwise evolution across virus species) were mapped onto published protein structures available for the S proteins of the four human-infecting betacoronaviruses studied here (Figure 5, Supplementary Figure 7). Finally, to compare dN/dS distributions between specific domains of the S protein within and across virus species, sites inferred to be under

15

1   positive or negative selection (PSS, NSS) were mapped onto S protein structures

2   (Supplementary Data 2).

3

4   **6. Validation through resampling and by comparing mutational distributions**

5   To validate our initial observations derived from virus genomes sampled up to May 2021, we

6   sought to determine if the 30 mutations that had been identified initially were also present in the

7   expanded embecov- and sarbecovirus diversity sampled up to July 25[th] 2022 (corresponding to

8   the final sampling date of this study). Virus diversity now included genome sequences derived

9   from more recently collected human isolates (only made publicly available after our initial

10  sampling), and from other closely related embeco- and sarbecoviruses from non-human hosts.

11  The expanded alignment comprises 1455 sequences (~700 embecovirus + SARS-CoV and

12  ~700 SARS-CoV-2), resulting in 27503 columns that were re-aligned under a progressive

13  profile-to-profile approach based on taxonomic relatedness to be further used to estimate an

14  expanded 'Maximum Likelihood' tree (following Methods section 2). To additionally explore if the

15  mutations identified were also present in a larger dataset representing an expanded SARS-

16  CoV-2 diversity (sampled up to July 25[th] 2022), a set of 1400 SARS-CoV-2 genomes denoting

17  'evolutionary successful' virus lineages (Supplementary Table 3) was examined independently

18  (Supplementary Text 3, Supplementary Figure 4 and 5). Both datasets were analysed following

19  the steps described in Methods Section 2 and 3, specifically searching for the mutations listed in

20  Table 1. Virus taxa included in both re-sampled datasets are listed in Supplementary Data 5. A

21  full description of the sequence subsampling and methodological approach used is available as

22  Supplementary Text 3, and Supplementary Figures 4 and 5.

23  We further sought to explore if the proportion of mutational patterns we classified as

24  putatively denoting homoplasy and/or stepwise evolution were more likely to arise from an

25  evolutionary scenario mostly driven by selection, compared to 'random' mutational patterns

26  derived from evolutionary scenarios generally driven by genetic drift. For this purpose, the

27  expanded alignment was translated to amino acid sequences and used to simulate three

28  alignments with 'AliSim' (http://www.iqtree.org/doc/AliSim) under the 'mimick real alignment'

29  function (mimicking a 'real' evolutionary process based on amino acid evolution under a LG

30  model, and applied to the inputted original tree). To compare the corresponding proportion of

31  sites scored under homoplasy and/or stepwise evolution, each dataset (the expanded and three

32  simulated alignments) was analysed following the steps described in Methods section 3. The

33  classification of mutational patterns within expanded and simulated datasets also serves the

34  purpose of validating our algorithm, originally developed for analysing the global dataset (that

16

1  included only OC43, HKU1, SARS-CoV-1 and SARS-CoV-2 sampled from the human host).

2  Associated results and a brief discussion are available as Supplementary Text 5.

3

### 4  7.  Reconstruction of amino acid evolution for selected sites

5  To further confirm our ML-derived results (see Methods section 3), for those mutations

6  displaying cumulative evidence of adaptive convergence (18121, 21623, 21635 and 23948,

7  Table 1), we used the expanded dataset to infer ancestral states under a Bayesian framework.

8  For each site, we first estimated a MCC (maximum clade credibility) tree from the resampled

9  codon alignment using an SRD06 substitution model (Shapiro et al. 2006) and a strict molecular

10 clock. Coded amino acid traits were then mapped onto the nodes of the MCC tree by performing

11 reconstructions of ancestral states under an asymmetric discrete trait evolution model (DTA) in

12 BEAST v1.8.4 (Lemey et al. 2009; Suchard et al. 2018). The DTA model was run using a

13 Bayesian Skygrid tree prior for $100 \times 10^6$ generations and sampled every 10,000 states until all

14 DTA-relevant parameters reached an ESS >200. For illustrative purposes, Figure 3 only shows

15 sites 18121, 21623 and 23948. The amino acid evolution pattern observed for site 21635 is

16 available in Supplementary Data 3.

17

17

1  **FIGURE LEGENDS**
2
3  **Fig 1. Phylogenetic tree of human-infecting betacoronaviruses**. The expanded tree estimated
4  from the Orf1ab+S alignment comprising 1455 sequences (see Methods section 6), summarizing the
5  phylogenetic pattern observed for four distantly related human-infecting betacoronaviruses: HKU1, OC43,
6  SARS-CoV-1 and SARS-CoV-2. MERS and related virus sequences were included in the tree for rooting
7  purposes only. Both the *Embecovirus* subgenus (HKU1 and OC43 and related viruses) and the
8  *Sarbecovirus* subgenus (SARS-CoV-1 and SARS-CoV-2 and related viruses) are indicated, showing the
9  positioning of the most closely related virus genome sequences derived from animal isolates (when
10 available). The different genotypes identified for the HKU1 (A, B and C) and for the OC43 (A–H) are
11 shown in Supplementary Data 1.
12
13 **Fig 2. Distribution of conserved/variable sites with S across different virus species** (a)
14 Top-down (upper panel) and side view (bottom panel) of a cartoon representation of the multidomain
15 architecture of the trimeric SARS-CoV-2 S ectodomain (PDB: 6VXX). The S1 subunit is coloured
16 according to the different protein domains: S1[A] in cream, S1[B] in teal, S1[C] in orange, and S1[D] in blue,
17 whilst the S2 subunit is shown in grey. (b) Top-down and side views of sphere-based representations of
18 trimeric S protein ectodomains for the viruses studied here: SARS-CoV-2 (PDB: 6VXX), SARS-CoV-1
19 (PDB: 6ACC), OC43 (PDB: 6OHW) and HKU1 (PDB: 5I08). The sphere-based representation shows
20 conserved (shown in grey; residues present ≥99% of all sequences) and variable sites (blue; residues
21 present in ≥1% of all sequences) across virus species. Variable sites identified as denoting homoplasy or
22 stepwise evolutionary patterns are shown in red (see Methods section 3). The asparagine residues of N-
23 linked glycosylation sequons are indicated in purple.
24
25 **Fig 3. Reconstruction of amino acid evolution at selected sites**. Maximum clade credibility
26 (MCC) trees showing the evolutionary trajectories for different amino acid states observed for three
27 illustrative sites (18121, 21623 and 23948) that (i) display evidence of homoplasy and/or stepwise
28 evolution, (ii) show evidence of positive selection, and (iii) are proximal to regions of established protein
29 function. The reconstructions of ancestral states for these sites show different amino acid states at nodes
30 (represented by circles in different colours). The posterior probability for a given amino acid state
31 occurring at a given node of interest is indicated. Sites 18121 display evidence of homoplasy across virus
32 lineages, site 21623 shows evidence of both homoplasy across species and stepwise evolution within
33 single virus species (*i.e.* OC43), and site 23948 shows evidence of stepwise evolution within single virus
34 species (*i.e.* SARS-CoV-1), and also of homoplasy across virus species (*i.e.* SARS-1 and SARS-CoV-2).
35
36 **Fig 4. Residue Ser[28] of nsp14 is situated near the nsp14-nsp10 interface.** Cartoon
37 representation of the SARS-CoV-1 nsp14-nsp10 protein complex (PDB: 5C8S) with Ser[28] (corresponding
38 to site 18121 in SARS-CoV-2 genome coordinates) shown as a red sphere. This residue is located within
39 the nsp14 ExoN domain (cream) and is approximately 9 Å from the interface with nsp10 (light blue, the
40 proximal nsp10 residue Cys[41] was used to calculate the distance and is indicated as a sphere). The
41 distance between nsp14's Ser[28] and nsp10's Cys[41] is annotated and indicated by a dashed black line.
42 Zoomed-in panel: detailed representation of the intra-nsp14 hydrogen-bond between the side chain of
43 Ser[28] and the main chain of Thr[25] (identified with the PISAebi server). The side chain of Ser[25] is indicated
44 as a red stick and Thr[25] is indicated in sticks and coloured according to atom (C, cream; O, red; N, blue).
45 The hydrogen-bond is indicated as a dashed black line.
46
47 **Fig 5. Mapping of mutations exhibiting homoplasy onto the S protein structure of SARS-**
48 **CoV-2.** Top-down (left) and side view (right) of a cartoon representation of the multidomain architecture
49 of the trimeric SARS-CoV-2 S ectodomain (PDB: 6VXX). The S2 subunit is highlighted in grey and the S1
50 ectodomain is divided into S1[A] (highlighted in cream), S1[B] (teal), S1[C] (orange), and S1[D] (blue) domains,
51 following the colour scheme in Figure 3. Homoplasic mutations that co-localize to known functional
52 surfaces (see Supplementary Table 4) are indicated in the structure and coloured in groups: Arg[21]
53 (corresponding to site 21623 in SARS-CoV-2 genome coordinates, in green), Pro[25] (site 21635, in green),
54 Asp[796] (site 23948, in yellow), Ile[1018] (site 24614, in red), Ala[1020] (site 24620, in red) and Leu[1024] (site
55 24632, in red). All representations are shown with a transparent protein surface for clarity.

18

## COMPETING INTERESTS

The authors declare no competing interests.

## DATA AVAILABILITY

Taxa IDs and accession numbers for virus sequences used in this study are provided in the Supplementary Data 4 and 5 files. All SARS-CoV-2 genome sequences and associated metadata used in this study are published in GISAID's EpiCoV database under the EPI SET GISAID Identifier: EPI_SET_230131zy. To view the contributors of each individual sequence with details such as accession number, virus name, collection date, originating and submitting lab, as well as the list of all authors, visit 10.55876/gis8.230131zy. PBD files used are listed as follows: S protein (HKU1 PDB:5I08, OC43 PDB:6OHW, SARS-CoV-1 PDB:6ACC and SARS-CoV-2 PDB:6VXX, 6ZGI). Orf1a (SARS-CoV-1 nsp3 PDB:2W2G). Orf1b (SARS-CoV-2 nsp13 PDB:6XEZ, SARS-CoV-1 nsp14 PDB:5C8S and SARS-CoV-2 nsp15 PDB:6WLC). Full code for our algorithm is available as open source: https://github.com/nataliamv/SARS-CoV-2-ARTs-Classification. An interactive notebook with our full selection analysis results is available at https://observablehq.com/@spond/beta-cov-analysis.

REFERENCES

Andersen KG, Rambaut A, Lipkin WI, Holmes EC, Garry RF. 2020. The proximal origin of SARS-CoV-2. Nat. Med. 26:450–452. doi: 10.1038/S41591-020-0820-9.

Avanzato VA et al. 2019. A structural basis for antibody-mediated neutralization of Nipah virus reveals a site of vulnerability at the fusion glycoprotein apex. Proc. Natl. Acad. Sci. U. S. A. 116:25057–25067. doi: 10.1073/PNAS.1912503116/-/DCSUPPLEMENTAL.

Banerjee A, Doxey AC, Mossman K, Irving AT. 2021. Unraveling the Zoonotic Origin and Transmission of SARS-CoV-2. Trends Ecol. Evol. 36:180–184. doi: 10.1016/J.TREE.2020.12.002.

Bedford T. 2021. Genetic diversity of betacoronaviruses including novel coronavirus (nCoV) 2020. https://nextstrain.org/groups/blab/beta-cov (Accessed December 23, 2021).

Boni MF et al. 2020. Evolutionary origins of the SARS-CoV-2 sarbecovirus lineage responsible for the COVID-19 pandemic. Nat. Microbiol. 5:1408–1417. doi: 10.1038/S41564-020-0771-4.

Campbell KM, Steiner G, Wells DK, Ribas A, Kalbasi A. 2020. Prioritization of SARS-CoV-2 epitopes using a pan-HLA and global population inference approach. bioRxiv. doi: 10.1101/2020.03.30.016931.

Cheng VCC, Lau SKP, Woo PCY, Kwok YY. 2007. Severe acute respiratory syndrome coronavirus as an agent of emerging and reemerging infection. Clin. Microbiol. Rev. 20:660–694. doi: 10.1128/CMR.00023-07.

COG-UK. 2021. COVID-19 Genomics UK Consortium. https://www.cogconsortium.uk/ (Accessed December 23, 2021).

Corman VM, Muth D, Niemeyer D, Drosten C. 2018. Hosts and Sources of Endemic Human Coronaviruses. Adv. Virus Res. 100:163–188. doi: 10.1016/BS.AIVIR.2018.01.001.

Delport W, Scheffler K, Seoighe C. 2008. Frequent toggling between alternative amino acids is driven by selection in HIV-1. PLoS Pathog. 4:e1000242. doi: 10.1371/JOURNAL.PPAT.1000242.

Dolan PT, Whitfield ZJ, Andino R. 2018. Mapping the Evolutionary Potential of RNA Viruses. Cell Host Microbe. 23:435–446. doi: 10.1016/J.CHOM.2018.03.012.

van Dorp L et al. 2020. No evidence for increased transmissibility from recurrent mutations in SARS-CoV-2. Nat. Commun. 11:5986. doi: 2.

Ellegren H. 2008. Comparative genomics and the study of evolution by natural selection. Mol. Ecol. 17:4586–4596. doi: 10.1111/J.1365-294X.2008.03954.X.

Escalera-Zamudio M et al. 2021. Identification of site-specific evolutionary trajectories shared across human betacoronaviruses. bioRxiv. doi: 10.1101/2021.05.24.445313.

Escalera-Zamudio M et al. 2020. Parallel evolution in the emergence of highly pathogenic avian influenza A viruses. Nat. Commun. 11:1–11. doi: 10.1038/s41467-020-19364-x.

Faria NR et al. 2021. Genomics and epidemiology of the P.1 SARS-CoV-2 lineage in Manaus, Brazil. Science. 372:815–821. doi: 10.1126/SCIENCE.ABH2644.

Fehr AR, Channappanavar R, Perlman S. 2017. Middle East Respiratory Syndrome: Emergence of a Pathogenic Human Coronavirus. Annu. Rev. Med. 68:387–399. doi: 10.1146/annurev-med-051215-031152.

Forni D et al. 2021. Adaptation of the endemic coronaviruses HCoV-OC43 and HCoV-229E to the human host. Virus Evol. 7:veab061. doi: 10.1093/VE/VEAB061.

Di Giorgio S, Martignano F, Torcia MG, Mattiuz G, Conticello SG. 2020. Evidence for host-dependent regular RNA editing in the transcriptome of SARS-CoV-2. Sci. Adv. 6:eabb5813. doi: 10.1007/s13353-022-00688-x.

GISAID. 2021. Global Initiative on Sharing Avian Influenza Data. https://www.gisaid.org/ (Accessed December 23, 2021).

Gutierrez B et al. 2022. Emergence and widespread circulation of a recombinant SARS-CoV-2 lineage in North America. Cell Host Microbe. 30:1112-1123.e3. doi: 10.1016/J.CHOM.2022.06.010.

Gutierrez B, Escalera-Zamudio M, Pybus OG. 2019. Parallel molecular evolution and adaptation in viruses. Curr. Opin. Virol. 34:90–96. doi: 10.1016/j.coviro.2018.12.006.

Hackbart M, Deng X, Baker SC. 2020. Coronavirus endoribonuclease targets viral polyuridine sequences to evade activating host sensors. Proc. Natl. Acad. Sci. U. S. A. 117:8094–8103. doi: 10.1073/PNAS.1921485117.

He JF et al. 2004. Molecular evolution of the SARS coronavirus during the course of the SARS epidemic in China. Science. 303:1666–1669. doi: 10.1126/SCIENCE.1092002.

Holmes EC, Rambaut A. 2004. Viral evolution and the emergence of SARS coronavirus. Philos. Trans. R. Soc. Lond. B. Biol. Sci. 359:1059–1065. doi: 10.1098/RSTB.2004.1478.

Hulswit RJG et al. 2019. Human coronaviruses OC43 and HKU1 bind to 9-O-acetylated sialic acids via a conserved receptor-binding site in spike protein domain A. Proc. Natl. Acad. Sci. U. S. A. 116:2681–2690. doi: 10.1073/PNAS.1809667116/-/DCSUPPLEMENTAL.

Hulswit RJG, De Haan CAM, Bosch B-J. 2016. Coronavirus Spike Protein and Tropism Changes. Adv. Virus Res. 96:29–57. doi: 10.1016/bs.aivir.2016.08.004.

HyPhy. 2013. PRIME. http://hyphy.org/w/index.php/PRIME (Accessed December 23, 2021).

ICTV et al. 2017. *2017.013S, Nidovirales*.

Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7: Improvements in performance and usability. Mol. Biol. Evol. 30:772–780. doi: 10.1093/molbev/mst010.

Kearse M et al. 2012. Geneious Basic: an integrated and extendable desktop software platform for the organization and analysis of sequence data. Bioinformatics. 28:1647–1649. doi: 10.1093/BIOINFORMATICS/BTS199.

Kemp SA et al. 2021. SARS-CoV-2 evolution during treatment of chronic infection. Nature. 592:277–282. doi: 10.1038/s41586-021-03291-y.

Kirchdoerfer RN et al. 2016. Pre-fusion structure of a human coronavirus spike protein. Nature. 531:118–121. doi: 10.1038/nature17200.

Kissler SM, Tedijanto C, Goldstein E, Grad YH, Lipsitch M. 2020. Projecting the transmission dynamics of SARS-CoV-2 through the postpandemic period. Science. 368:860–868. doi: 10.1126/science.abb5793.

Kistler KE, Bedford T. 2021. Evidence for adaptive evolution in the receptor-binding domain of seasonal coronaviruses OC43 and 229E. Elife. 10:1–35. doi: 10.7554/eLife.64509.

1 Kosakovsky Pond SL, Frost SDW, et al. 2006. Adaptation to different human populations by HIV-1 revealed by codon-based
2 analyses. PLoS Comput. Biol. 2:e62. doi: 10.1371/JOURNAL.PCBI.0020062.
3 Kosakovsky Pond SL. Evolutionary annotation of global SARS-CoV-2/COVID-19 genomes enabled by data from GSAID 2021.
4 https://observablehq.com/@spond/sars_cov_2_sites (Accessed December 23, 2021).
5 Kosakovsky Pond SL, Frost SDW. 2005. Not so different after all: a comparison of methods for detecting amino acid sites under
6 selection. Mol. Biol. Evol. 22:1208–1222. doi: 10.1093/MOLBEV/MSI105.
7 Kosakovsky Pond SL, Murrell B, Poon AFY. 2012. Evolution of viral genomes: interplay between selection, recombination, and other
8 forces. Methods Mol. Biol. 856:239–272. doi: 10.1007/978-1-61779-585-5_10.
9 Kosakovsky Pond SL, Posada D, Gravenor MB, Woelk CH, Frost SDW. 2006. GARD: a genetic algorithm for recombination
10 detection. Bioinformatics. 22:3096–3098. doi: 10.1093/BIOINFORMATICS/BTL474.
11 Kosakovsky Pond SL, Wisotsky SR, Escalante A, Magalis BR, Weaver S. 2021. Contrast-FEL-A Test for Differences in Selective
12 Pressures at Individual Sites among Clades and Sets of Branches. Mol. Biol. Evol. 38:1184–1198. doi:
13 10.1093/MOLBEV/MSAA263.
14 Krissinel E, Henrick K. 2007. Inference of macromolecular assemblies from crystalline state. J. Mol. Biol. 372:774–797. doi:
15 10.1016/J.JMB.2007.05.022.
16 Lan J et al. 2020. Structure of the SARS-CoV-2 spike receptor-binding domain bound to the ACE2 receptor. Nature. 581:215–220.
17 doi: 10.1038/s41586-020-2180-5.
18 Lau SKP et al. 2015. Discovery of a Novel Coronavirus, China Rattus Coronavirus HKU24, from Norway Rats Supports the Murine
19 Origin of Betacoronavirus 1 and Has Implications for the Ancestor of Betacoronavirus Lineage A. J. Virol. 89:3076–3092. doi:
20 10.1128/JVI.02420-14.
21 Lau SKP et al. 2011. Molecular epidemiology of human coronavirus OC43 reveals evolution of different genotypes over time and
22 recent emergence of a novel genotype due to natural recombination. J. Virol. 85:11325–11337. doi: 10.1128/JVI.05512-11.
23 Lemey P, Rambaut A, Drummond AJ, Suchard MA. 2009. Bayesian phylogeography finds its roots. PLoS Comput. Biol.
24 5:e1000520. doi: 10.1371/JOURNAL.PCBI.1000520.
25 Li F. 2016. Structure, Function, and Evolution of Coronavirus Spike Proteins. Annu. Rev. Virol. 3:237–261. doi: 10.1146/annurev-
26 virology-110615-042301.
27 Li Fang, Li W, Farzan M, Harrison SC. 2005. Structure of SARS coronavirus spike receptor-binding domain complexed with
28 receptor. Science. 309:1864–1868. doi: 10.1126/SCIENCE.1116480.
29 Li Wendong et al. 2005. Bats are natural reservoirs of SARS-like coronaviruses. Science. 310:676–679. doi:
30 10.1126/SCIENCE.1118391.
31 Loewe L, Hill WG. 2010. The population genetics of mutations: good, bad and indifferent. Philos. Trans. R. Soc. Lond. B. Biol. Sci.
32 365:1153–1167. doi: 10.1098/RSTB.2009.0317.
33 Ma Y et al. 2015. Structural basis and functional analysis of the SARS coronavirus nsp14-nsp10 complex. Proc. Natl. Acad. Sci. U.
34 S. A. 112:9436–9441. doi: 10.1073/PNAS.1508686112.
35 MacLean OA et al. 2021. Natural selection in the evolution of SARS-CoV-2 in bats created a generalist virus and highly capable
36 human pathogen. PLoS Biol. 19:e3001115. doi: 10.1371/JOURNAL.PBIO.3001115.
37 De Maio N et al. 2021. Mutation Rates and Selection on Synonymous Mutations in SARS-CoV-2. Genome Biol. Evol. 13:evab087.
38 doi: 10.1093/GBE/EVAB087.
39 McIntosh K, Becker WB, Chanock RM. 1967. Growth in suckling-mouse brain of 'IBV-like' viruses from patients with upper
40 respiratory tract disease. Proc. Natl. Acad. Sci. U. S. A. 58:2268–2273. doi: 10.1073/PNAS.58.6.2268.
41 Menachery VD, Graham RL, Baric RS. 2017. Jumping species - a mechanism for coronavirus persistence and survival. Curr. Opin.
42 Virol. 23:1–7. doi: 10.1016/j.coviro.2017.01.002.
43 Millet JK, Whittaker GR. 2015. Host cell proteases: Critical determinants of coronavirus tropism and pathogenesis. Virus Res.
44 202:120–134.
45 Murrell B et al. 2012. Detecting individual sites subject to episodic diversifying selection. PLoS Genet. 8:e1002764. doi:
46 10.1371/JOURNAL.PGEN.1002764.
47 Murrell B et al. 2015. Gene-Wide Identification of Episodic Selection. Mol. Biol. Evol. 32:1365–1371. doi:
48 10.1093/MOLBEV/MSV035.
49 Nelde A et al. 2021. SARS-CoV-2-derived peptides define heterologous and COVID-19-induced T cell recognition. Nat. Immunol.
50 22:74–85. doi: 10.1038/S41590-020-00808-X.
51 O'Toole Á et al. 2021. Tracking the international spread of SARS-CoV-2 lineages B.1.1.7 and B.1.351/501Y-V2 with grinch.
52 Wellcome Open Res. 6:121. doi: 10.12688/WELLCOMEOPENRES.16661.2.
53 Ohta T. 1973. Slightly deleterious mutant substitutions in evolution. Nature. 246:96–98. doi: 10.1038/246096A0.
54 Oong XY et al. 2017. Identification and evolutionary dynamics of two novel human coronavirus OC43 genotypes associated with
55 acute respiratory infections: phylogenetic, spatiotemporal and transmission network analyses. Emerg. Microbes Infect. 6:e3. doi:
56 10.1038/EMI.2016.132.
57 Peiris JSM et al. 2003. Coronavirus as a possible cause of severe acute respiratory syndrome. Lancet. 361:1319–1325. doi:
58 10.1016/S0140-6736(03)13077-2.
59 Pollett S et al. 2021. A comparative recombination analysis of human coronaviruses and implications for the SARS-CoV-2
60 pandemic. Sci. Rep. 11:17365. doi: 10.1038/s41598-021-96626-8.
61 Rausch JW, Capoferri AA, Katusiime MG, Patro SC, Kearney MF. 2020. Low genetic diversity may be an Achilles heel of SARS-
62 CoV-2. Proc. Natl. Acad. Sci. U. S. A. 117:24614–24616. doi: 10.1073/PNAS.2017726117.
63 Rochman ND et al. 2022. Epistasis at the SARS-CoV-2 Receptor-Binding Domain Interface and the Propitiously Boring Implications
64 for Vaccine Escape. MBio. 13:e00135-22. doi: 10.1128/mbio.00135-22.
65 Sagulenko P, Puller V, Neher RA. 2018. TreeTime: Maximum-likelihood phylodynamic analysis. Virus Evol. 4:vex042. doi:
66 10.1093/ve/vex042.
67 Shaman J, Galanti M. 2020. Will SARS-CoV-2 become endemic? Science. 370:527–529. doi: 10.1126/SCIENCE.ABE5960.
68 Shang J et al. 2020. Structural basis of receptor recognition by SARS-CoV-2. Nature. 581:221–224. doi: 10.1038/s41586-020-2179-
69 y.
70 Shapiro B, Rambaut A, Drummond AJ. 2006. Choosing appropriate substitution models for the phylogenetic analysis of protein-

21

coding sequences. Mol. Biol. Evol. 23:7–9. doi: 10.1093/MOLBEV/MSJ021.

Simmonds P. 2020. Rampant C→U Hypermutation in the Genomes of SARS-CoV-2 and Other Coronaviruses: Causes and Consequences for Their Short- and Long-Term Evolutionary Trajectories. mSphere. 5:e00408–e00420. doi: 10.1128/MSPHERE.00408-20.

Stamatakis A. 2015. Using RAxML to Infer Phylogenies. Curr. Protoc. Bioinforma. 51:6.14.1-6.14.14. doi: 10.1002/0471250953.BI0614S51.

Starr TN et al. 2022. Shifting mutational constraints in the SARS-CoV-2 receptor-binding domain during viral evolution. Science. 377:420–424. doi: 10.1126/SCIENCE.ABO7896/SUPPL_FILE/SCIENCE.ABO7896_DATA_S1.ZIP.

Stern A et al. 2017. The Evolutionary Pathway to Virulence of an RNA Virus. Cell. 169:35-46.e19. doi: 10.1016/J.CELL.2017.03.013.

Su S et al. 2016. Epidemiology, Genetic Recombination, and Pathogenesis of Coronaviruses. Trends Microbiol. 24:490–502. doi: 10.1016/j.tim.2016.03.003.

Suchard MA et al. 2018. Bayesian phylogenetic and phylodynamic data integration using BEAST 1.10. Virus Evol. 4:vey016. doi: 10.1093/VE/VEY016.

Taefehshokr N, Taefehshokr S, Hemmat N, Heit B. 2020. Covid-19: Perspectives on Innate Immune Evasion. Front. Immunol. 11:580641. doi: 10.3389/FIMMU.2020.580641.

Tegally H et al. 2021. Detection of a SARS-CoV-2 variant of concern in South Africa. Nature. 592:438–443. doi: 10.1038/S41586-021-03402-9.

Tortorici MA et al. 2019. Structural basis for human coronavirus attachment to sialic acid receptors. Nat. Struct. Mol. Biol. 26:481–489. doi: 10.1038/s41594-019-0233-y.

Turakhia Y et al. 2022. Pandemic-scale phylogenomics reveals the SARS-CoV-2 recombination landscape. Nature. 609:994–997. doi: 10.1038/s41586-022-05189-9.

V'kovski P, Kratzel A, Steiner S, Stalder H, Thiel V. 2020. Coronavirus biology and replication: implications for SARS-CoV-2. Nat. Rev. Microbiol. 19:155–170. doi: 10.1038/s41579-020-00468-6.

Vijaykrishna D et al. 2007. Evolutionary Insights into the Ecology of Coronaviruses. J. Virol. 81:4012–4020. doi: 10.1128/jvi.02605-06.

Vijgen L et al. 2005. Complete Genomic Sequence of Human Coronavirus OC43: Molecular Clock Analysis Suggests a Relatively Recent Zoonotic Coronavirus Transmission Event. J. Virol. 79:1595–1604. doi: 10.1128/jvi.79.3.1595-1604.2005.

ViPR-NCBI. 2021. Virus Pathogen Resource. https://www.viprbrc.org/brc/home.spg?decorator=vipr (Accessed December 23, 2021).

Wang G, Dunbrack Jr RL. 2004. Scoring profile-to-profile sequence alignments. Protein Sci. 13:1612–1626. doi: 10.1110/PS.03601504.

Wang H, Pipes L, Nielsen R. 2020. Synonymous mutations and the molecular evolution of SARS-CoV-2 origins. Virus Evol. 7:veaa098. doi: 10.1093/VE/VEAA098.

Wang Y et al. 2015. Coronavirus nsp10/nsp16 Methyltransferase Can Be Targeted by nsp10-Derived Peptide In Vitro and In Vivo To Reduce Replication and Pathogenesis. J. Virol. 89:8416–8427. doi: 10.1128/JVI.00948-15.

Watanabe Y, Bowden TA, Wilson IA, Crispin M. 2019. Exploitation of glycosylation in enveloped virus pathobiology. Biochim. Biophys. Acta Gen. Subj. 1863:1480–1497. doi: 10.1016/j.bbagen.2019.05.012.

WHO. 2021. Middle East respiratory syndrome 2021. http://www.emro.who.int/health-topics/mers-cov/mers-outbreaks.html (Accessed December 23, 2021).

de Wilde AH, Snijder EJ, Kikkert M, van Hemert MJ. 2018. Host Factors in Coronavirus Replication. Curr. Top. Microbiol. Immunol. 419:1–42. doi: 10.1007/82_2017_25.

Woo PCY et al. 2005. Characterization and Complete Genome Sequence of a Novel Coronavirus, Coronavirus HKU1, from Patients with Pneumonia. J. Virol. 79:884–895. doi: 10.1128/jvi.79.2.884-895.2005.

Woo PCY et al. 2006. Comparative Analysis of 22 Coronavirus HKU1 Genomes Reveals a Novel Genotype and Evidence of Natural Recombination in Coronavirus HKU1. J. Virol. 80:7136–7145. doi: 10.1128/jvi.00509-06.

Woo PCY, Huang Y, Lau SKP, Yuen K-Y. 2010. Coronavirus genomics and bioinformatics analysis. Viruses. 2:1804–1820. doi: 10.3390/V2081803.

Yewdell JW. 2021. Antigenic drift: Understanding COVID-19. Immunity. 54:2681–2687. doi: 10.1016/j.immuni.2021.11.016.

Yuen CK et al. 2020. SARS-CoV-2 nsp13, nsp14, nsp15 and orf6 function as potent interferon antagonists. Emerg. Microbes Infect. 9:1418–1428. doi: 10.1080/22221751.2020.1780953.

Zahradník J et al. 2021. SARS-CoV-2 variant prediction and antiviral drug design are enabled by RBD in vitro evolution. Nat. Microbiol. 6:1188–1198. doi: 10.1038/s41564-021-00954-4.

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

Zhu Y et al. 2018. A novel human coronavirus OC43 genotype detected in mainland China. Emerg. Microbes Infect. 7:173. doi: 10.1038/S41426-018-0171-5.

22

## Table 1. Potentially relevant sites across human-infecting betacoronaviruses

| SARS-CoV-2 genome coordinates † | ORF | Protein/ Residue † | Mutations observed (global tree) | | | | | | Confirmed in re-sampled trees: | | Homoplasy (H)/ Stepwise Evolution (SWE) | Selection across species (Method, p-value) †# | Selection in SARS-CoV-2 (recent amino acid changes)¶ | Epitopes |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Ancestral LinA | OC43 | HKU1 | Ancestral LinB | SARS-CoV-1 | SARS-CoV-2§ | Expanded | 1400 SARS-CoV-2 | | | | |
| 2557 | Orf1a | nsp2 585 | P | P | S | S | A/T | P/S | Yes | Conserved in SARS-CoV-2 | H/SWE | | | 0 |
| 7478 | Orf1a | nsp3 1587 | N | S/N | N | N | T | N | Yes | Conserved in SARS-CoV-2 | H | | PSS, N→S/D (OC43-like and new state) | 0 |
| 16189 | Orf1b | nsp12 917 | D | D | E/D | E | E | E | Yes | Conserved in SARS-CoV-2 | H/SWE | Overall negative selection (FEL 0.02) | | 1 |
| 17809 | Orf1b | nsp13 525 | V | V | V/I | I | I | I | Yes | Conserved in SARS-CoV-2 | H | | | 0 |
| **18121** | **Orf1b** | **nsp14 28** | **A** | **A** | **A/S** | **S** | **S** | **S** | Yes | Conserved in SARS-CoV-2 | **H** | **Different overall positive selection (CF 0.022)** | | **1** |
| 18334 | Orf1b | nsp14 100 | D | D | E/D | E | D | E | Yes | Conserved in SARS-CoV-2 | H/SWE | Overall negative selection (FEL 0.004) | | 0 |
| 18442 | Orf1b | nsp14 136 | K | K | K/R | R | R | R | Yes | Conserved in SARS-CoV-2 | H | | | 0 |
| 19048 | Orf1b | nsp14 338 | A | G/A | G | A | A | A | Yes | Conserved in SARS-CoV-2 | H/SWE | OC43 branch (MEME 0.035) | | 0 |
| 20344 | Orf1b | nsp15 243 | Q | Q | H/Y | H | H | H | Yes | Conserved in SARS-CoV-2 | H/SWE | | | 2 |
| 20554 | Orf1b | nsp15 313 | N | N | S/N | S | S | S | No | Conserved in SARS-CoV-2 | H | Overall negative selection (FEL 0.04) | | 0 |
| 21400 | Orf1b | nsp16 249 | A | A | T/S | S | S | S | Yes | Conserved in SARS-CoV-2 | H/SWE | | | 2 |
| 21614 | Orf S | S1 18 | F | F/I/L | I | L | F | L | Yes | Variable in SARS-CoV-2 | H/SWE | | PSS, L→F (OC43 and SARS-CoV-1-like) | 1 |
| **21623** | **Orf S** | **S1 21** | **V** | **R/V/K /I** | **K/Y/L** | **R** | **V** | **R/I** | Yes | Conserved in SARS-CoV-2 | **H/SWE** | **HKU1, OC43 and SARS-2 branches (MEME 0.047)** | **NSS, R→ I/K/T (OC43 and HKU1-like and new state)** | **1** |
| **21635** | **Orf S** | **S1 25** | **V** | **P/V/S /L/H** | **V/I** | **P** | **N** | **P/S** | Yes | Conserved in SARS-CoV-2 | **H/SWE** | **HKU1, OC43 and SARS-2 branches (MEME 0.048)** | **NSS, P→S and L (OC43-like)** | 0 |
| 21800 | Orf S | S1 81 | K | K | Q/K | D | G/D | D | Yes | Variable in SARS-CoV-2 | SWE | | PSS, D→Y/A/G (SARS-CoV-1-like and new states) | 0 |
| 21863 | Orf S | S1 102 | Y | F/I/T | Y | I | V | I | Yes | Conserved in SARS-CoV-2 | H/SWE | | PSS, I→V (SARS-CoV-1-like) | 0 |
| 21920 | Orf S | S1 120 | V | V | V/I | V | I | V | Yes | Conserved in SARS-CoV-2 | H/SWE | | | 0 |
| 21926 | Orf S | S1 122 | T | T | N/T | N | N | N | Yes | Conserved in SARS-CoV-2 | H/SWE | Overall negative selection (FEL 0.002) | NSS | 0 |
| 22004 | Orf S | S1 149 | N | N/K | K/I | N | G | N | Yes | Conserved in SARS-CoV-2 | H | | NSS, N→D (new state) | 0 |
| 22124 | Orf S | S1 189 | D | T/D/N | D | H | H | N | Yes | Conserved in SARS-CoV-2 | H/SWE | OC43 branch (MEME 0.008) | NSS | 0 |
| 22553 | Orf S | S1 332 | N | D/N | D/N | N | N | N | Yes | Conserved in SARS-CoV-2 | H/SWE | | | 1 |
| 23048 | Orf S | S1 497 | S | A/G/S | D/S | G | G | G | Yes | Variable in SARS-CoV-2 | H/SWE | HKU1 branch (MEME 0.044) | | 2 |
| **23948** | **Orf S** | **S2 796** | **D** | **N** | **D** | **D** | **Y/D** | **Y/D** | D/N discrepancy | Variable in SARS-CoV-2 | **SWE** | **Different overall positive selection (CF 0.031)** | **PSS, D → Y/G/H (SARS-CoV-1-like and new states)** | **0** |
| 24614 | Orf S | S2 1018 | V | V | V/I | I | I | I | Yes | Conserved in SARS-CoV-2 | H | | NSS | 1 |
| 24620 | Orf S | S2 1020 | F | F | F/A/L | A | A | A | Yes | Conserved in SARS-CoV-2 | H | | PSS, A→S/V (new states) | 2 |
| 24632 | Orf S | S2 1024 | Q | Q | L/R | L | L | L | Yes | Conserved in SARS-CoV-2 | H/SWE | | | 2 |
| 24863 | Orf S | S2 1101 | T | T | H/S | H | S | H | Yes | Conserved in SARS-CoV-2 | H/SWE | | NSS, H→Y (new state) | 1 |
| 25037 | Orf S | S2 1159 | Q | Q | Q/H | H | H | H | Yes | Conserved in SARS-CoV-2 | H | | NSS, H→Y (new state) | 0 |
| 25166 | Orf S | S2 1202 | D | D/Y | D/E | E | E | E | Yes | Conserved in SARS-CoV-2 | H | | PSS, E→Q/G (new states) | 0 |
| 25247 | Orf S | S2 1230 | V | V | V/M | M | M | M | Yes | Conserved in SARS-CoV-2 | H | | PSS, M→I/T/ L (new states) | 1 |

23

† Positions indicate the start of the codon for reference genome Wuhan-Hu-1 (NC_045512.2). Sites in bold refer to those highlighted in the results section
# Sites/branches scored under MEME/FEL and Contrast-FEL (CF); CF tests for differences is selective pressures between clades
§ Representing virus diversity sampled as of May 2021
¶ Representing viral diversity sampled as of December 2022 available from: https://observablehq.com/@spond/sars_cov_2_sites
* Potential T cell epitopes derived from HLA class I and HLA-DR SARS-CoV-2 binding peptides (Campbell et al. 2020; Nelde et al. 2021)
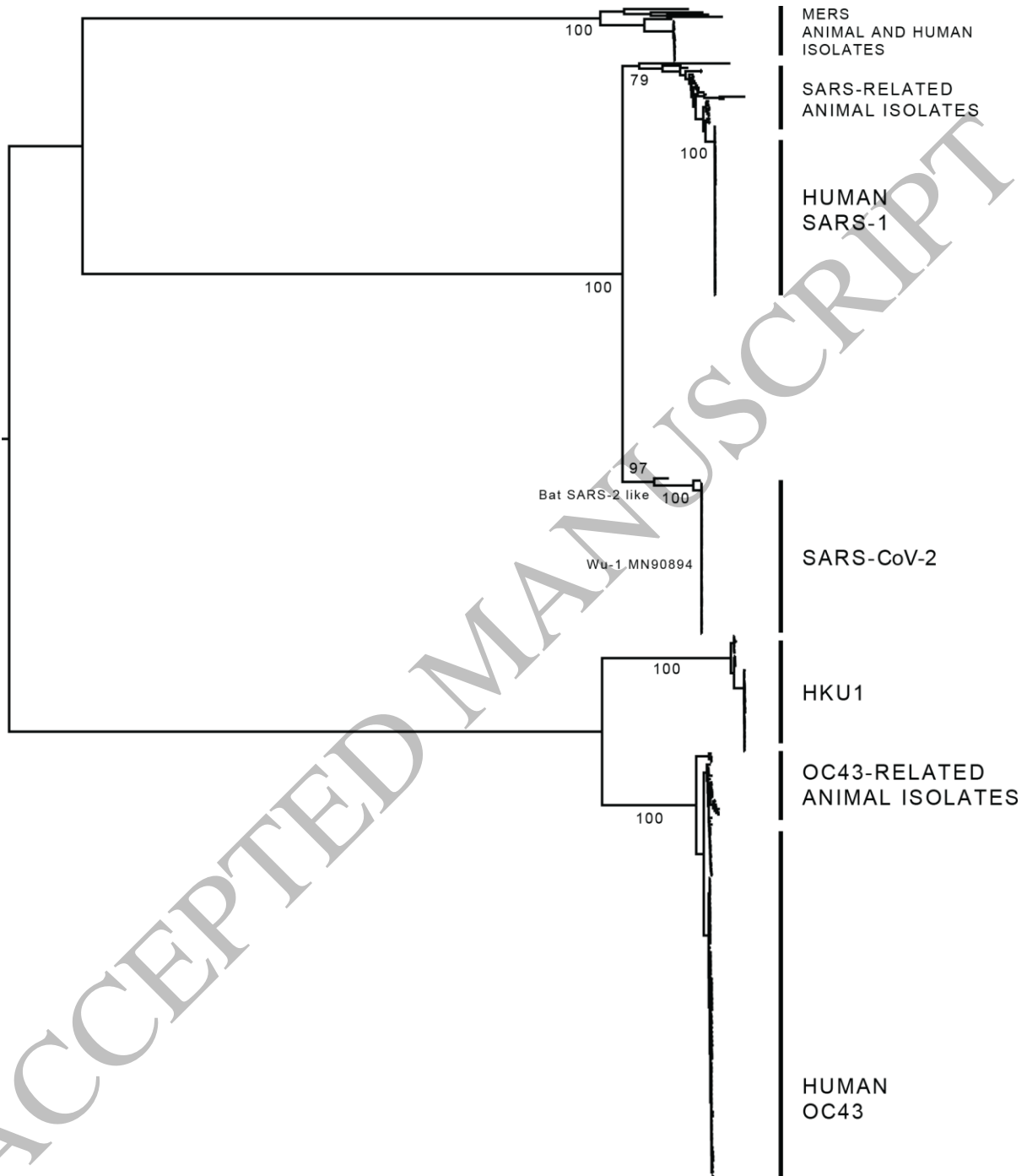
1



MERS
ANIMAL AND HUMAN
ISOLATES

SARS-RELATED
ANIMAL ISOLATES

HUMAN
SARS-1

Bat SARS-2 like

Wu-1 MN90894

SARS-CoV-2

HKU1

OC43-RELATED
ANIMAL ISOLATES

HUMAN
OC43

0.5

2

3
4

*Figure 1*
*159x190 mm ( x  DPI)*

5

**A**

Top View

S1^A
S1^C
S1^B

Side View

S1^B
S1^A
S1^C
S1^D

S1
S2

■ Conserved (within virus species)
■ Variable (within virus species)
■ Sites evidencing putative homoplasy/stepwise evolution across virus species
■ N-glycan Asn (within virus species)

**B**

SARS-CoV-2    SARS-CoV-1    HCoV-OC43    HCoV-HKU1

*Figure 2*
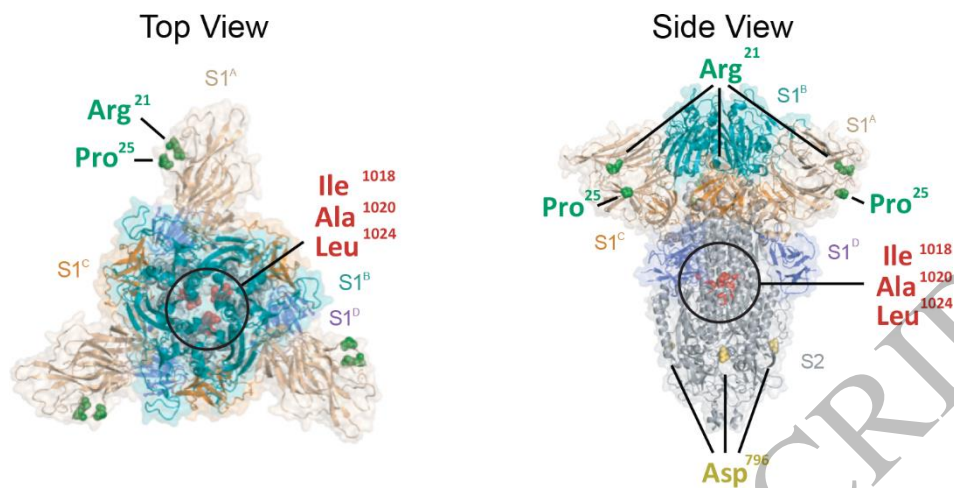*159x74 mm ( x  DPI)*

1

2

3

4

*Figure 3*
*159x54 mm ( x  DPI)*

*Figure 4*
*77x65 mm ( x  DPI)*

1
2
3

4

*Figure 5*
*125x62 mm ( x  DPI)*