

 Open access • Posted Content • DOI:10.1101/2020.01.24.919183

Complete genome characterisation of a novel coronavirus associated with severe human respiratory disease in Wuhan, China — [Source link](#)

[Fan Wu](#), [Su Zhao](#), [Bin Yu](#), [Yan-Mei Chen](#) ...+16 more authors

Institutions: [Fudan University](#), [Huazhong University of Science and Technology](#), [Centers for Disease Control and Prevention](#), [Chinese Center for Disease Control and Prevention](#) ...+1 more institutions

Published on: 01 Jan 2020 - [bioRxiv](#) (Cold Spring Harbor Laboratory)

Topics: [Coronavirus](#) and [Outbreak](#)

Related papers:

- [Clinical features of patients infected with 2019 novel coronavirus in Wuhan, China](#)
- [Genomic characterisation and epidemiology of 2019 novel coronavirus: implications for virus origins and receptor binding.](#)
- [A Novel Coronavirus from Patients with Pneumonia in China, 2019.](#)
- [A pneumonia outbreak associated with a new coronavirus of probable bat origin](#)
- [Epidemiological and clinical characteristics of 99 cases of 2019 novel coronavirus pneumonia in Wuhan, China: a descriptive study](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/complete-genome-characterisation-of-a-novel-coronavirus-2h8f9fmpif>

1 Note: This paper has been revised after peer review, so that it can be considered technically
2 correct.

3

4 **Complete genome characterisation of a novel coronavirus**
5 **associated with severe human respiratory disease in Wuhan,**
6 **China**

7 Fan Wu^{1,5}, Su Zhao^{2,5}, Bin Yu^{3,5}, Yan-Mei Chen^{1,5}, Wen Wang^{1,5}, Yi Hu^{2,5}, Zhi-Gang Song^{1,5},
8 Zhao-Wu Tao², Jun-Hua Tian³, Yuan-Yuan Pei¹, Ming-Li Yuan², Yu-Ling Zhang¹, Fa-Hui
9 Dai¹, Yi Liu¹, Qi-Min Wang¹, Jiao-Jiao Zheng¹, Lin Xu¹, Edward C. Holmes⁴, Yong-Zhen
10 Zhang^{1*}

11

12 ¹Shanghai Public Health Clinical Center & School of Public Health, Fudan University,
13 Shanghai, China.

14 ²Department of Pulmonary and Critical Care Medicine, The Central Hospital of Wuhan,
15 Tongji Medical College, Huazhong University of Science and Technology, Wuhan 430014,
16 China.

17 ³Wuhan Center for Disease Control and Prevention, Wuhan, Hubei, China

18 ⁴Marie Bashir Institute for Infectious Diseases and Biosecurity, School of Life and
19 Environmental Sciences and School of Medical Sciences, The University of Sydney, Sydney,
20 Australia.

21

22 ⁵These authors contributed equally: Fan Wu, Su Zhao, Bin Yu, Yan-Mei Chen, Wen Wang, Yi
23 Hu, Zhi-Gang Song. *e-mail: zhangyongzhen@shphc.org.cn

24

25 **Emerging and re-emerging infectious diseases, such as SARS, MERS, Zika and highly**
26 **pathogenic influenza present a major threat to public health¹⁻³. Despite intense research**
27 **effort, how, when and where novel diseases appear are still the source of considerable**
28 **uncertainty. A severe respiratory disease was recently reported in the city of Wuhan,**
29 **Hubei province, China. At the time of writing, at least 62 suspected cases have been**
30 **reported since the first patient was hospitalized on December 12nd 2019. Epidemiological**
31 **investigation by the local Center for Disease Control and Prevention (CDC) suggested**
32 **that the outbreak was associated with a sea food market in Wuhan. We studied seven**
33 **patients who were workers at the market, and collected bronchoalveolar lavage fluid**
34 **(BALF) from one patient who exhibited a severe respiratory syndrome including fever,**
35 **dizziness and cough, and who was admitted to Wuhan Central Hospital on December**
36 **26th 2019. Next generation metagenomic RNA sequencing⁴ identified a novel RNA virus**
37 **from the family *Coronaviridae* designed WH-Human-1 coronavirus (WHCV).**
38 **Phylogenetic analysis of the complete viral genome (29,903 nucleotides) revealed that**
39 **WHCV was most closely related (89.1% nucleotide similarity similarity) to a group of**
40 **Severe Acute Respiratory Syndrome (SARS)-like coronaviruses (genus *Betacoronavirus*,**
41 **subgenus *Sarbecovirus*) previously sampled from bats in China and that have a history**
42 **of genomic recombination. This outbreak highlights the ongoing capacity of viral spill-**
43 **over from animals to cause severe disease in humans.**

44
45 Seven patients, comprising five men and two women, were hospitalized at the Central
46 Hospital of Wuhan from December 14 through December 28, 2019. The median age of the

47 patients was 43, ranging from 31 to 70 years old. The clinical characteristics of the patients
48 are shown in Table 1. Fever and cough were the most common symptoms. All patients had
49 fever with body temperatures ranging from 37.2°C to 40°C. Patients 1, 2, 5, 6 and 7 had
50 cough, while patients 1, 2 and 7 presented with severe cough with phlegm at onset of illness.
51 Patients 4 and 5 also complained of chest tightness and dyspnea. Patients 1, 3, 4 and 6
52 experienced dizziness and patient 3 felt weakness. No neurological symptoms were observed
53 in any of the patients. Bacterial culture revealed the presence of *Streptococcus* bacteria in
54 throat swabs from patients 3, 4 and 7. Combination antibiotic, antiviral and glucocorticoid
55 therapy were administered. Unfortunately, patient 1 and 4 showed respiratory failure: patient
56 1 was given high flow noninvasive ventilation, while patient 4 was provided with nasal/face
57 mask ventilation (Table 1).

58 Epidemiological investigation by the Wuhan CDC revealed that all the suspected cases
59 were linked to individuals working in a local indoor seafood market. Notably, in addition to
60 fish and shell fish, a variety of live wild animals including hedgehogs, badgers, snakes, and
61 birds (turtledoves) were available for sale in the market before the outbreak began, as well as
62 animal carcasses and animal meat. No bats were available for sale. While the patients might
63 have had contact with wild animals in the market, none recalled exposure to live poultry.

64 Patient 1 was a 41-year-old man with no history of hepatitis, tuberculosis or diabetes.
65 He was admitted and hospitalized in Wuhan Central Hospital 6 days after the onset of illness.
66 The patient reported fever, chest tightness, unproductive cough, pain and weakness for one
67 week on presentation. Physical examination of cardiovascular, abdominal and neurologic
68 examination was normal. Mild lymphopenia (less than 900 cells per cubic milli-meter) was

69 observed, but white blood cell and blood platelet count was normal in a complete blood count
70 (CBC) test. Elevated levels of C-reactive protein (CRP, 41.4 mg/L of blood, reference range
71 0-6 mg/L) was observed and levels of aspartate aminotransferase, lactic dehydrogenase, and
72 creatine kinase were slightly elevated in blood chemistry tests. The patient had mild
73 hypoxemia with oxygen levels of 67mmHg by the Arterial Blood Gas (ABG) Test. On the
74 first day of admission (day 6 after the onset of illness), chest radiographs were abnormal with
75 air-space shadowing such a ground-glass opacities, focal consolidation and patchy
76 consolidation in both lungs (Figure 1). Chest computed tomographic (CT) scans revealed
77 bilateral focal consolidation, lobar consolidation and patchy consolidation, especially in the
78 lower lung. A chest radiograph revealed a bilateral diffuse patchy and fuzzy shadow on day 5
79 after admission (day 11 after the onset of illness). Preliminary aetiological investigation
80 excluded the presence of influenza virus, *Chlamydia pneumoniae* and *Mycoplasma*
81 *pneumoniae* by commercial pathogen antigen detection kits and confirmed by PCR. Other
82 common respiratory pathogens, including adenovirus, were also negative by qPCR (Figure
83 S1). The condition of the patient did not improve after three days of treatment with combined
84 antiviral and antibiotic therapy. He was admitted to the intensive care unit (ICU) and
85 treatment with a high flow non-invasive ventilator was initiated. The patient was transferred
86 to another hospital in Wuhan for further treatment 6 days after admission.

87 To investigate the possible aetiologic agents associated this disease, we collected
88 bronchoalveolar lavage fluid (BALF) from patient 1 and performed deep meta-transcriptomic
89 sequencing. All the clinical specimens were handled in a biosafety level 3 laboratory at the
90 Shanghai Public Health Clinical Center. Total RNA was extracted from 200µl BAL fluid and

91 a meta-transcriptomic library was constructed for pair-end (150 bp) sequencing using an
92 Illumina MiniSeq as previously described⁴⁻⁷. In total, we generated 56,565,928 sequence
93 reads that were *de novo* assembled and screened for potential aetiologic agents. Of the
94 384,096 contigs assembled by Megahit⁸, the longest (30,474 nucleotides [nt]) had high
95 abundance and was closely related to a bat SARS-like coronavirus isolate - bat-SL-CoVZC45
96 (GenBank Accession MG772933) - previously sampled in China, with a nt identity of 89.1%
97 (Table S1 and S2). The genome sequence of this novel virus, as well as its termini, were
98 determined and confirmed by RT-PCR⁹ and 5'/3' RACE kits (TaKaRa), respectively. This
99 new virus was designated as WH-Human 1 coronavirus (WHCV) (and has also been referred
100 to as '2019-nCoV') and its whole genome sequence (29,903 nt) has been assigned GenBank
101 accession number MN908947. Remapping the RNA-seq data against the complete genome of
102 WHCV resulted in an assembly of 123,613 reads, providing 99.99% genome coverage at a
103 mean depth of 6.04X (range: 0.01X -78.84X) (Figure S2). The viral load in the BALF sample
104 was estimated by quantitative PCR (qPCR) to be 3.95×10^8 copies/mL (Figure S3).

105 The viral genome organization of WHCV was characterized by sequence alignment
106 against two representative members of the genus *Betacoronavirus*: a human-origin
107 coronavirus (SARS-CoV Tor2, AY274119) and a bat-origin coronavirus (Bat-SL-CoVZC45,
108 MG772933) (Figure 2). The un-translational regions (UTR) and open reading frame (ORF) of
109 WHCV were mapped based on this sequence alignment and ORF prediction. The WHCV
110 viral genome was similar to these two coronaviruses (Figure 2 and Table S3), with a gene
111 order 5'-replicase ORF1ab-S-envelope(E)-membrane(M)-N-3'. WHCV has 5' and 3' terminal
112 sequences typical of the betacoronaviruses, with 265 nt at the 5' terminal and 229 nt at the 3'

113 terminal region. The predicted replicase ORF1ab gene of WHCV is 21,291 nt in length and
114 contained 16 predicted non-structural proteins (Table S4), followed by (at least) 13
115 downstream ORFs. Additionally, WHCV shares a highly conserved domain (LLRKNGNKG:
116 amino acids 122-130) with SARS-CoV in nsp1. The predicted S, ORF3a, E, M and N genes
117 of WHCV are 3,822, 828, 228, 669 and 1,260 nt in length, respectively. In addition to these
118 ORFs regions that are shared by all members of the subgenus *Sarbecovirus*, WHCV is similar
119 to SARS-CoV in that it carries a predicted ORF8 gene (366 nt in length) located between the
120 M and N ORF genes. The functions of WHCV ORFs were predicted based on those of known
121 coronaviruses and given in Table S5. In a manner similar to SARS CoV Tor2, a leader
122 transcription regulatory sequence (TRS) and nine putative body TRSs could be readily
123 identified upstream of the 5' end of ORF, with the putative conserved TRS core sequence
124 appeared in two forms – the ACGAAC or CUAAAC (Table S6).

125 To determine the evolutionary relationships between WHCV and previously identified
126 coronaviruses, we estimated phylogenetic trees based on the nucleotide sequences of the
127 whole genome sequence, non-structural protein genes ORF1a and 1b, and the main structural
128 proteins encoded by the S, E, M and N genes (Figures 3 and S4). In all phylogenies WHCV
129 clustered with members of the subgenus *Sarbecovirus*, including the SARS-CoV responsible
130 for the global SARS pandemic of 2002-2003^{1,2}, as well as a number of SARS-like
131 coronaviruses sampled from bats. However, WHCV changed topological position within the
132 subgenus *Sarbecovirus* depending on which gene was used, suggestive of a past history of
133 recombination in this group of viruses (Figures 3 and S4). Specifically, in the S gene tree
134 (Figure S4), WHCV was most closely related to the bat coronavirus bat-SL-CoVZC45 with

135 82.3% amino acid (aa) identity (and ~77.2% aa identity to SARS CoV; Table S3), while in the
136 ORF1b phylogeny WHCV fell in a basal position within the subgenus *Sarbecovirus* (Figure
137 3). This topological division was also observed in the phylogenetic trees estimated for
138 conserved domains in the replicase polyprotein pp1ab (Figure S5).

139 To better understand the potential of WHCV to infect humans, the receptor-binding
140 domain (RBD) of its spike protein was compared to those in SARS-CoVs and bat SARS-like
141 CoVs. The RBD sequences of WHCV were more closely related to those of SARS-CoVs
142 (73.8%-74.9% aa identity) and SARS-like CoVs including strains Rs4874, Rs7327 and
143 Rs4231 (75.9%-76.9% aa identity) that are able to use the human ACE2 receptor for cell entry
144 (Table S7)¹⁰. In addition, the WHCV RBD was only one amino acid longer than the SARS-
145 CoV RBD (Figure 4a). In contrast, other bat SARS-like CoVs including the Rp3 strain that
146 cannot use human ACE2¹¹, had amino acid deletions at positions 473-477 and 460-472
147 compared to the SARS-CoVs (Figure 4a). The previously determined¹² crystal structure of
148 SARS-CoV RBD complexed with human ACE2 (PDB 2AJF) revealed that regions 473-477
149 and 460-472 directly interact with human ACE2 and hence may be important in determining
150 species specificity (Figure 4b). We predicted the three-dimension protein structures of WHCV,
151 Rs4874 and Rp3 RBD domains by protein homology modelling using the SWISS-MODEL
152 server and compared them to the crystal structure of SARS-CoV RBD domains (PDB 2GHV)
153 (Figure 4, c-f). In accord with the sequence alignment, the predicted protein structures of
154 WHCV and Rs4874 RBD domains were closely related to that of SARS-CoVs and different
155 from the predicted structure of the RBD domain from Rp3. In addition, the N-terminus of
156 WHCV S protein is more similar to that of SARS-CoV rather than other human coronaviruses

157 (HKU1 and OC43) (Figure S6) that can bind to sialic acid¹³. In sum, the high similarities of
158 amino acid sequences and predicted protein structure between WHCV and SARS-CoV RBD
159 domains suggest that WHCV may efficiently use human ACE2 as a cellular entry receptor,
160 perhaps facilitating human-to-human transmission^{10, 14-15}.

161 To further characterize putative recombination events in the evolutionary history of the
162 sarbecoviruses the whole genome sequence of WHCV and four representative coronaviruses -
163 Bat SARS-like CoV Rp3, CoVZC45, CoVZXC21 and SARS-CoV Tor2 - were analysed using
164 the Recombination Detection Program v4 (RDP4)¹⁶. Although the similarity plots suggested
165 possible recombination events between WHCV and SARS CoVs or SARS-like CoVs (Figure
166 S7), there was no significant evidence for recombination across the genome as a whole.
167 However, some evidence for past recombination was detected in the S gene of WHCV and
168 SARS CoV and bat SARS-like CoVs (WIV1 and RsSHC014) ($p < 3.147 \times 10^{-3}$ to $p < 9.198 \times 10^{-9}$),
169 with similarity plots suggesting the presence of recombination break points at nucleotides
170 1,029 and 1,652 that separated the WHCV S gene into three regions (Figure 5). In
171 phylogenies of the fragment nt 1 to 1029 and nt 1652 to the end of the sequence, WHCV was
172 most closely related to Bat-SL-CoVZC45 and Bat-SL-CoVZXC21, whereas in the region nt
173 1030 to 1651 (the RBD region) WHCV grouped with SARS CoV and bat SARS-like CoVs
174 (WIV1 and RsSHC014) that are capable of direct human transmission^{14,17}.

175 Coronaviruses are associated with a number of infectious disease outbreaks in humans,
176 including SARS in 2002/3 and MERS in 2012^{1,18}. Four other coronaviruses - human
177 coronaviruses HKU1, OC43, NL63 and 229E - are also associated with respiratory disease¹⁹⁻
178 ²². Although SARS-like coronaviruses have been widely identified in mammals including bats

179 since 2005 in China^{9,23-25}, the exact origin of human-infected coronaviruses remains unclear.
180 Herein, we describe a novel coronavirus - WHCV - in BALF from a patient experiencing
181 severe respiratory disease in Wuhan, China. Phylogenetic analysis suggested that WHCV
182 represents a novel virus within genus *Betacoronavirus* (subgenus *Sarbecovirus*) and hence
183 that exhibits some genomic and phylogenetic similarity to SARS-CoV¹, particularly in the
184 RBD. These genomic and clinical similarities to SARS, as well as its high abundance in
185 clinical samples, provides evidence for an association between WHCV and the ongoing
186 outbreak of respiratory disease in Wuhan.

187 The identification of multiple SARS-like-CoVs in bats led to the idea that these animals
188 act as the natural reservoir hosts of these viruses^{19,20}. Although SARS-like viruses have been
189 identified widely in bats in China, viruses identical to SARS-CoV have not yet been
190 documented. Notably, WHCV is most closely related to bat coronaviruses, even exhibiting
191 100% aa similarity to Bat-SL-CoVZC45 in the nsp7 and E proteins. Hence, these data suggest
192 that bats are a possible reservoir host of WHCV. However, as a variety of animal species were
193 for sale in the market when the disease was first reported, more work is needed to determine
194 the natural reservoir and any intermediate hosts of WHCV.

195

196 **Acknowledgements** This study was supported by the Special National Project on
197 investigation of basic resources of China (Grant SQ2019FY010009) and the National Natural
198 Science Foundation of China (Grants 81861138003 and 31930001). ECH is supported by an
199 ARC Australian Laureate Fellowship (FL170100022).

200

201 **Author Contributions** Y.-Z.Z. conceived and designed the study. S.Z, Y.H, Z.-W.T. and M.-
202 L.Y. performed the clinical work and sample collection. B.Y and J.-H.T. performed
203 epidemiological investigation and sample collection. F.W, Z.-G.S., L.X., Y.-Y.P., Y.-L.Z., F.-
204 H.D., Y.L., J.-J.Z. and Q.-M.W. performed the experiments. Y.-M.C., W.W., F.W., E.C.H. and
205 Y.-Z.Z. analysed the data. Y.-Z.Z. E.C.H. and F.W. wrote the paper with input from all
206 authors. Y.-Z.Z. led the study.

207

208 Correspondence and requests for materials should be addressed to Y.-Z.Z.

209 (zhangyongzhen@shphc.org.cn).

210

211 **Competing interests** The authors declare no competing interests.

212 **References**

- 213 1. Drosten, C., et al., Identification of a novel coronavirus in patients with severe acute
214 respiratory syndrome. *N. Engl. J. Med.* **348**,1967–1976 (2003).
- 215 2. Wolfe, N.D., Dunavan, C.P., Diamond, J. Origins of major human infectious diseases.
216 *Nature.* **447**, 279-283 (2007).
- 217 3. Ventura, C.V., Maia, M., Bravo-Filho, V., Góis, A.L. & Belfort, R. Jr. Zika virus in Brazil
218 and macular atrophy in a child with microcephaly. *Lancet.* **387**, 228 (2016).
- 219 4. Shi, M. et al. Redefining the invertebrate RNA virosphere. *Nature.* **540**, 539-543 (2016).
- 220 5. Shi, M., et al. The evolutionary history of vertebrate RNA viruses. *Nature.* **556**,197-202
221 (2018).
- 222 6. Yadav, P.D., et al. Nipah virus sequences from humans and bats during Nipah outbreak,
223 Kerala, India, 2018. *Emerg. Infect. Dis.* **25**, 1003-1006 (2019).
- 224 7. McMullan, L.K., et al. Characterisation of infectious Ebola virus from the ongoing
225 outbreak to guide response activities in the Democratic Republic of the Congo: a
226 phylogenetic and in vitro analysis. *Lancet. Infect. Dis.* **19**, 1023-1032 (2019).
- 227 8. Li, D., Liu, C.M., Luo, R., Sadakane, K. & Lam, T.W. MEGAHIT: An ultra-fast single-
228 node solution for large and complex metagenomics assembly via succinct de Bruijn
229 graph. *Bioinformatics* **31**, 1674-1676 (2015).
- 230 9. Wang, W., et al. Discovery, diversity and evolution of novel coronaviruses sampled from
231 rodents in China. *Virology.* **474**, 19-27 (2015).
- 232 10. Hu, B. et al. Discovery of a rich gene pool of bat SARS-related coronaviruses provides
233 new insights into the origin of SARS coronavirus. *PLoS Pathog.* **13**: e1006698 (2017).

- 234 11. Ren, W. et al. Difference in receptor usage between severe acute respiratory syndrome
235 (SARS) coronavirus and SARS-like coronavirus of bat origin. *J Virol.* **82**:1899-1907
236 (2008).
- 237 12. Li, F., Li, W., Farzan, M., Harrison, S.C. Structure of SARS coronavirus spike receptor-
238 binding domain complexed with receptor. *Science.* **309**, 1864-1868 (2005).
- 239 13. Hulswit, R.J.G., et al. Human coronaviruses OC43 and HKU1 bind to 9-O-acetylated
240 sialic acids via a conserved receptor-binding site in spike protein domain A. *Proc Natl*
241 *Acad Sci USA.*, **116**, 2681-2690 (2019).
- 242 14. Ge, X.Y. et al. Isolation and characterization of a bat SARS-like coronavirus that uses the
243 ACE2 receptor. *Nature.* **503**: 535-538 (2013).
- 244 15. Yang, X.L., et al. Isolation and characterization of a novel bat coronavirus closely related
245 to the direct progenitor of severe acute respiratory syndrome coronavirus. *J Virol.* **90**:
246 3253-3256 (2016).
- 247 16. Martin, D.P., Lemey, P., Lott, M., Moulton, V., Posada, D., Lefevre, P. RDP3: a flexible
248 and fast computer program for analyzing recombination. *Bioinformatics* 26:2462–2463
249 (2010).
- 250 17. Menachery, V.D., et al. A SARS-like cluster of circulating bat coronaviruses shows
251 potential for human emergence. *Nat Med.* **21**:1508-1513 (2015).
- 252 18. Bermingham, A., et al. Severe respiratory illness caused by a novel coronavirus, in a
253 patient transferred to the United Kingdom from the Middle East, September 2012. *Euro.*
254 *Surveill.* **17**, 20290 (2012).
- 255 19. Hamre, D. & Procknow, J.J. A new virus isolated from the human respiratory tract. *Proc.*

- 256 *Soc. Exp. Biol. Med.* **121**, 190–193 (1966).
- 257 20. McIntosh, K., Becker, W.B., Chanock, R.M. Growth in suckling-mouse brain of "IBV-
258 like" viruses from patients with upper respiratory tract disease. *Proc Natl Acad Sci USA*.
259 **58**, 2268-73(1967).
- 260 21. van der Hoek, L., et al. Identification of a new human coronavirus. *Nat. Med.* **10**, 368–373
261 (2004).
- 262 22. Woo, P.C., et al. Characterization and complete genome sequence of a novel coronavirus,
263 coronavirus HKU1, from patients with pneumonia. *J. Virol.* **79**,884–895 (2005).
- 264 23. Li, W., et al. Bats are natural reservoirs of SARS-like coronaviruses. *Science* **310**, 676–
265 679 (2005).
- 266 24. Lau S.K., et al. Severe acute respiratory syndrome coronavirus- like virus in Chinese
267 horseshoe bats. *Proc.Natl.Acad.Sci.U.S.A.* **102**, 14040–14045 (2005).
- 268 25. Wang, W., et al. Discovery of a highly divergent coronavirus in the Asian house shrew
269 from China illuminates the origin of the Alphacoronaviruses. *J. Virol.* **91**, e00764-17
270 (2017).
- 271

272 **Figure legends**

273

274 **Figure 1.** Chest radiographs of patient 1. **a. b. c. d.** Chest computed tomographic scans of
275 Patient 1 were obtained on the day of admission (day 6 after the onset of illness). Bilateral
276 focal consolidation, lobar consolidation, and patchy consolidation were clearly observed,
277 especially in the lower lung. **e.** Chest radiograph of patient 1 was obtained on day 5 after
278 admission (day 11 after the onset of illness). Bilateral diffuse patchy and fuzzy shadow were
279 observed.

280

281 **Figure 2.** Genome organization of SARS and SARS-like CoVs including Tor2, CoVZC45
282 and WHCV determined here.

283

284 **Figure 3.** Maximum likelihood phylogenetic trees of nucleotide sequences of the ORF1a,
285 ORF1b, E and M genes of WHCV and related coronaviruses. Numbers (>70) above or below
286 branches indicate percentage bootstrap values for the associated nodes. The trees were mid-
287 point rooted for clarity only. The scale bar represents the number of substitutions per site.

288

289 **Figure 4.** Analysis of receptor-binding domain (RBD) of the spike (S) protein of WHCV
290 coronavirus. **(a)** Amino acid sequence alignment of SARS-like CoV RBD sequences. Three
291 bat SARS-like CoVs, which could efficiently utilize the human ACE2 as receptor, had an
292 RBD sequence of similar size to SARS-CoV, and WHCV contains a single Val 470 insertion.
293 The key amino acid residues involved in the interaction with human ACE2 are marked with a

294 brown box. In contrast, five bat SARS-like CoVs had amino acid deletions at two motifs
295 (amino acids 473-477 and 460-472) compared with those of SARS-CoV, and Rp3 has been
296 reported not to use ACE2.¹¹ **(b)** The two motifs (aa 473-477 and aa 460-472) are shown in red
297 on the crystal structure of the SARS-CoV spike RBD complexed with receptor human ACE2
298 (PDB 2AJF). Human ACE2 is shown in blue and the SARS-CoV spike RBD is shown in
299 green. Important residues in human ACE2 that interact with SARS-CoV spike RBD are
300 marked. **(c)** Predicted protein structures of RBD of WHCV spike protein based on target-
301 template alignment using ProMod3 on the SWISS-MODEL server. The most reliable models
302 were selected based on GMQE and QMEAN Scores. Template: 2ghw.1.A, GMQE: 0.83;
303 QMEAN:-2.67. Motifs resembling amino acids 473-477 and 460-472 of the SARS-CoV spike
304 protein are shown in red. **(d)** Predicted structure of RBD of SARS-like CoV Rs4874.
305 Template: 2ghw.1.A, GMQE:0.99; QMEAN:-0.72. Motifs resembling amino acids 473-477
306 and 460-472 of the SARS-CoV spike protein are shown in red. **(e)** Predicted structure of the
307 RBD of SARS-like CoV Rp3. Template: 2ghw.1.A, GMQE:0.81, QMEAN:-1.50. **(f)** Crystal
308 structure of RBD of SARS-CoV spike protein (green) (PDB 2GHV). Motifs of amino acids
309 473-477 and 460-472 are shown in red.

310

311 **Figure 5.** Possible recombination events in the S gene of sarbecoviruses. A sequence
312 similarity plot (upper panel) reveals two putative recombination break-points shown by black
313 dashed lines, with their locations indicated at the bottom. The plot shows S gene similarity
314 comparisons of the WHCV (query) against SARS-CoV Tor2 and bat SARS-like CoVs WIV1,
315 Rf1 and CoVZC45. Phylogenies of major parental region (1-1028 and 1653-3804) and minor

316 parental region (1029-1652) are shown below the similarity plot. Phylogenies were estimated
317 using a ML method and mid-point rooted for clarity only. Numbers above or below branches
318 indicate percentage bootstrap values.

319

320 **Methods**

321 **Cases and collection of clinical data and samples**

322 Patients presenting with acute onset of fever ($>37.5^{\circ}\text{C}$), cough, and chest tightness, and who
323 were admitted to Wuhan Central Hospital in Wuhan city, China, were considered as suspected
324 cases. During admission, bronchoalveolar lavage fluid (BALF) was collected and stored at -
325 80°C until further processing. Demographic, clinical and laboratory data were retrieved from
326 the clinical records of the confirmed patients.

327

328 **RNA library construction and sequencing**

329 Total RNA was extracted from the BALF sample of patient 1 using the RNeasy Plus
330 Universal Mini Kit (Qiagen) following the manufacturer's instructions. The quantity and
331 quality of the RNA solution was assessed using a Qbit machine and an Agilent 2100
332 Bioanalyzer (Agilent Technologies) before library construction and sequencing. An RNA
333 library was then constructed using the SMARTer Stranded Total RNA-Seq Kit v2 (TaKaRa,
334 Dalian, China). Ribosomal RNA (rRNA) depletion was performed during library construction
335 following the manufacturer's instructions. Paired-end (150 bp) sequencing of the RNA library
336 was performed on the MiniSeq platform (Illumina). Library preparation and sequencing were
337 carried out at the Shanghai Public Health Clinical Center, Fudan University, Shanghai, China.

338 **Data processing and viral agent identification**

339 Sequencing reads were first adaptor- and quality-trimmed using the Trimmomatic program²⁶.
340 The remaining reads (56,565,928 reads) were assembled *de novo* using both the Megahit
341 (version 1.1.3)⁸ and Trinity program (version 2.5.1)²⁷ with default parameter settings. Megahit
342 generated a total of 384,096 assembled contigs (size range: 200-30,474 nt), while Trinity
343 generated 1,329,960 contigs with a size range of 201 to 11,760 nt. All of these assembled
344 contigs were compared (using blastn and Diamond blastx) against the entire non-redundant
345 nucleotide (Nt) and protein (Nr) database, with *e*-values set to 1×10^{-10} and 1×10^{-5} ,
346 respectively. To identify possible aetiologic agents present in the sequence data, the
347 abundance of the assembled contigs was first evaluated as the expected counts using the
348 RSEM program²⁸ implemented in Trinity. Non-human reads (23,712,657 reads), generated by
349 filtering host reads using the human genome (human release 32, GRCh38.p13, downloaded
350 from Gencode) by Bowtie2²⁹, were used for the RSEM abundance assessment.

351 As the longest contigs generated by Megahit (30,474 nt) and Trinity (11,760 nt) both
352 had high similarity to the bat SARS-like coronavirus isolate bat-SL-CoVZC45 and were at
353 high abundance (Table S1 and S2), the longer one (30,474 nt) that covered almost the whole
354 virus genome was used for primer design for PCR confirmation and genome termini
355 determination. Primers used in PCR, qPCR and RACE experiments are listed in Table S8. The
356 PCR assay was conducted as described previously⁹ and the complete genome termini was
357 determined using the Takara SMARTer RACE 5'/3' kit (TaKaRa) following the
358 manufacturer's instructions. Subsequently, the genome coverage and sequencing depth were
359 determined by remapping all of the adaptor- and quality-trimmed reads to the whole genome

360 of WHCV using Bowtie2²⁹ and Samtools³⁰.

361 The viral loads of WHCV in BALF of patient 1 were determined by quantitative real-
362 time RT-PCR with Takara One Step PrimeScript™ RT-PCR Kit (Takara RR064A) following
363 the manufacturer's instructions. Real-time RT-PCR was performed using 2.5µl RNA with
364 8pmol of each primer and 4pmol probe under the following conditions: reverse transcription
365 at 42°C for 10 minutes, and 95°C for 1 minute, followed by 40 cycles of 95°C for 15 seconds
366 and 60°C for 1 minute. The reactions were performed and detected by ABI 7500 Real-Time
367 PCR Systems. PCR product covering the Taqman primers and probe region was cloned into
368 pLB vector using the Lethal Based Simple Fast Cloning Kit (TIAGEN) as standards for
369 quantitative viral load test.

370 **Virus genome characterization and phylogenetic analysis**

371 For the newly identified virus genome, the potential open reading frames (ORFs) were
372 predicted and annotated using the conserved signatures of the cleavage sites recognized by
373 coronavirus proteinases, and were processed in the Lasergene software package (version 7.1,
374 DNASTar). The viral genes were aligned using the L-INS-i algorithm implemented in MAFFT
375 (version 7.407)³¹.

376 Phylogenetic analyses were then performed using the nucleotide sequences of various
377 CoV gene data sets: (i) Whole genome, (ii) ORF1a, (iii) ORF1b, (iv) nsp5 (3CLpro), (v)
378 RdRp (nsp12), (vi) nsp13 (Hel), (vii) nsp14 (ExoN), (viii) nsp15 (NendoU), (ix) nsp16 (O-
379 MT), (x) spike (S), and the (xi) nucleocapsid (N). Phylogenetic trees were inferred using the
380 Maximum likelihood (ML) method implemented in the PhyML program (version 3.0)³², using

381 the Generalised Time Reversible substitution (GTR) model and Subtree Pruning and
382 Regrafting (SPR) branch-swapping. Bootstrap support values were calculated from 1,000
383 pseudo-replicate trees. The best-fit model of nucleotide substitution was determined using
384 MEGA (version 5)³³. Amino acid identities among sequences were calculated using the
385 MegAlign program implemented in the Lasergene software package (version 7.1, DNASTar).

386 **Genome recombination analysis**

387 Potential recombination events in the history of the sarbecoviruses were assessed using both
388 the Recombination Detection Program v4 (RDP4)¹⁶ and Simplot (version 3.5.1)³⁴. The RDP4
389 analysis was conducted based on the complete genome (nucleotide) sequence, employing the
390 RDP, GENECONV, BootScan, maximum chi square, Chimera, SISCAN, and 3SEQ methods.
391 Putative recombination events were identified with a Bonferroni corrected p-value cut-off of
392 0.01. Similarity plots were inferred using Simplot to further characterize potential
393 recombination events, including the location of breakpoints.

394 **Analysis of RBD domain of WHCV spike protein**

395 An amino acid sequence alignment of WHCV, SARS-CoVs, bat SARS-like CoVs RBD
396 sequences was performed using MUSCLE³⁵. The predicted protein structures of the spike
397 protein RBD were estimated based on target-template alignment using ProMod3 on SWISS-
398 MODEL server (<https://swissmodel.expasy.org/>). The sequences of the spike RBD domains of
399 WHCV, Rs4874 and Rp3 were searched by BLAST against the primary amino acid sequence
400 contained in the SWISS-MODEL template library (SMTL, last update: 2020-01-09, last
401 included PDB release: 2020-01-03). Models were built based on the target-template alignment
402 using ProMod3. The global and per-residue model quality were assessed using the QMEAN

403 scoring function³⁶. The PDB files of the predicted protein structures were displayed and
404 compared with the crystal structures of SARS-CoV spike RBD (PDB 2GHV)³⁷ and the crystal
405 of structure of SARS-CoV spike RBD complexed with human ACE2 (PDB 2AJF)¹².
406

407 **References**

- 408 26. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina
409 sequence data. *Bioinformatics* **30**, 2114-2120 (2014).
- 410 27. Grabherr, M. G. et al. Full-length transcriptome assembly from RNA-Seq data without a
411 reference genome. *Nat. Biotechnol.* **29**, 644–652 (2011).
- 412 28. Li, B., Ruotti, V., Stewart, R.M., Thomson, J.A. & Dewey, C.N. RNA-Seq gene
413 expression estimation with read mapping uncertainty. *Bioinformatics* **26**, 493-500 (2010).
- 414 29. Langmead, B. & Salzberg, S.L. Fast gapped-read alignment with Bowtie 2. *Nat Methods*.
415 **9**, 357–359 (2012).
- 416 30. Li, H. et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **15**,
417 2978-2079 (2009).
- 418 31. Katoh, K. & Standley, D.M. MAFFT multiple sequence alignment software version 7:
419 improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772-780 (2013).
- 420 32. Guindon, S. et al. New algorithms and methods to estimate maximum-likelihood
421 phylogenies: assessing the performance of PhyML 3.0. *Syst.Biol.* **59**, 307-321 (2010).
- 422 33. Tamura, K. et al. MEGA5: molecular evolutionary genetics analysis using maximum
423 likelihood, evolutionary distance, and maximum parsimony methods. *Mol. Biol. Evol.* **28**,
424 2731–2739 (2011).
- 425 34. Lole, K.S. et al. Full-length human immunodeficiency virus type 1 genomes from
426 subtype C-infected seroconverters in India, with evidence of intersubtype recombination.
427 *J. Virol.* **73**, 152–160 (1999).
- 428 35. Edgar, R.C. MUSCLE: multiple sequence alignment with high accuracy and high

- 429 throughput. *Nucleic Acids Res.* **32**, 1792-1797 (2004).
- 430 36. Waterhouse, A. et al. SWISS-MODEL: homology modelling of protein structures and
431 complexes. *Nucleic Acids Res.* **46**, W296-W303 (2018).
- 432 37. Hwang, W.C. et al. Structural basis of neutralization by a human anti-severe acute
433 respiratory syndrome spike protein antibody, 80R. *J Biol Chem.* **281**, 34610-34616
434 (2006).
- 435

436 **Supplementary legends**

437

438 **Supplementary Tables**

439 **Table S1.** The top 50 abundant assembled contigs generated using the Megahit program.

440 **Table S2.** The top 80 abundant assembled contigs generated using the Trinity program.

441 **Table S3.** Amino acid identities of the selected predicted gene products between the novel
442 coronavirus (WHCV) and known betacoronaviruses.

443 **Table S4.** Cleavage products of the replicase polyproteins of WHCV.

444 **Table S5.** Predicted gene functions of WHCV ORFs.

445 **Table S6.** Coding of potential and putative transcription regulatory sequences of the genome
446 sequence of WHCV.

447 **Table S7.** Amino acid identities of the RBD sequence between SARS- and bat SARS-like
448 CoVs.

449 **Table S8.** PCR primers used in this study.

450

451 **Supplementary Figures**

452 **Figure S1.** Detection of other respiratory pathogens by qPCR.

453 **Figure S2.** Mapped read count plot showing the coverage depth per base of the WHCV
454 genome.

455 **Figure S3.** Detection of WHCV in clinical samples by RT-qPCR. **(a)** Specificity of the
456 WHCV primers used in RT-qPCR. Test samples comprised clinical samples that are positive
457 for at least one of the following viruses: Influenza A virus (09H1N1 and H3N2), Influenza B
458 virus, Human adenovirus, Respiratory syncytial virus, Rhinovirus, Parainfluenza virus type 1-
459 4, Human bocavirus, Human metapneumovirus, Coronavirus OC43, Coronavirus NL63,
460 Coronavirus 229E and Coronavirus HKU1. **(b-c)** Standard curve. **(d)** Amplification curve of
461 WHCV.

462 **Figure S4.** Maximum likelihood phylogenetic trees of the nucleotide sequences of the whole
463 genome, S and N genes of WHCV and related coronaviruses. Numbers (>70) above or below
464 branches indicate percentage bootstrap values. The trees were mid-point rooted for clarity
465 only. The scale bar represents the number of substitutions per site.

466 **Figure S5.** Maximum likelihood phylogenetic trees of the nucleotide sequences of the 3CL,
467 RdRp, Hel, ExoN, NendoU, and O-MT genes of WHCV and related coronaviruses. Numbers
468 (>70) above or below branches indicate percentage bootstrap values. The trees were mid-point
469 rooted for clarity only. The scale bar represents the number of substitutions per site.

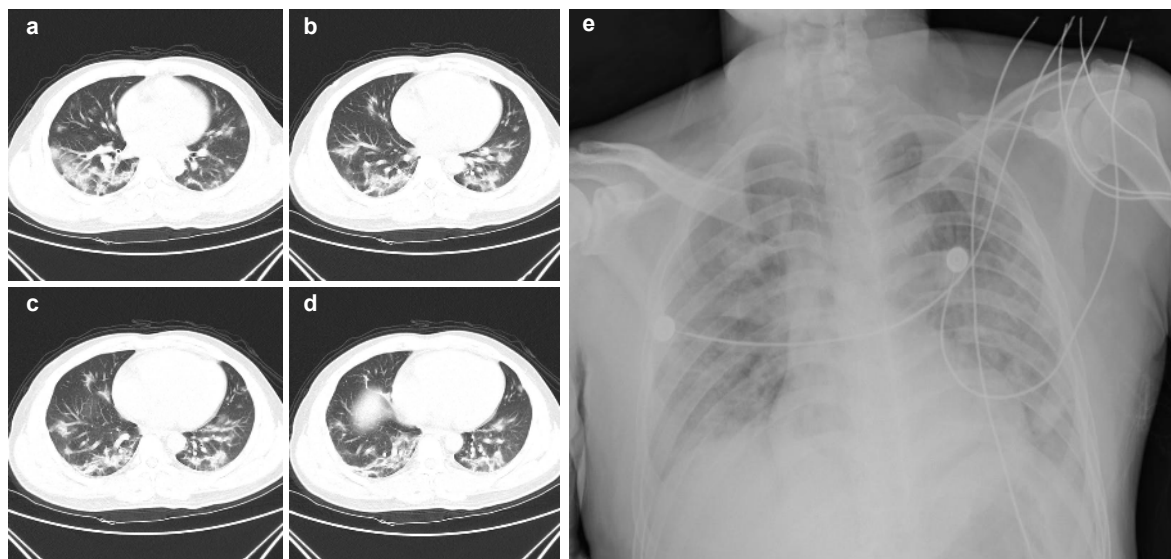
470 **Figure S6.** Amino acid sequence comparison of the N-terminal domain (NTD) of spike
471 protein of WHCV, bovine coronavirus (BCoV), mouse hepatitis virus (MHV) and human

472 coronavirus (HCoV OC43 and HKU1) that can bind to sialic acid and the SARS-CoVs that
473 cannot. The key residues¹³ for sialic acid binding on BCoV, MHV, HCoV OC43 and HKU1
474 were marked with a brown box.

475 **Figure S7.** A sequence similarity plot of WHCV, SARS- and bat SARS-like CoVs revealing
476 putative recombination events.

Table 1. Clinical symptoms and patient data

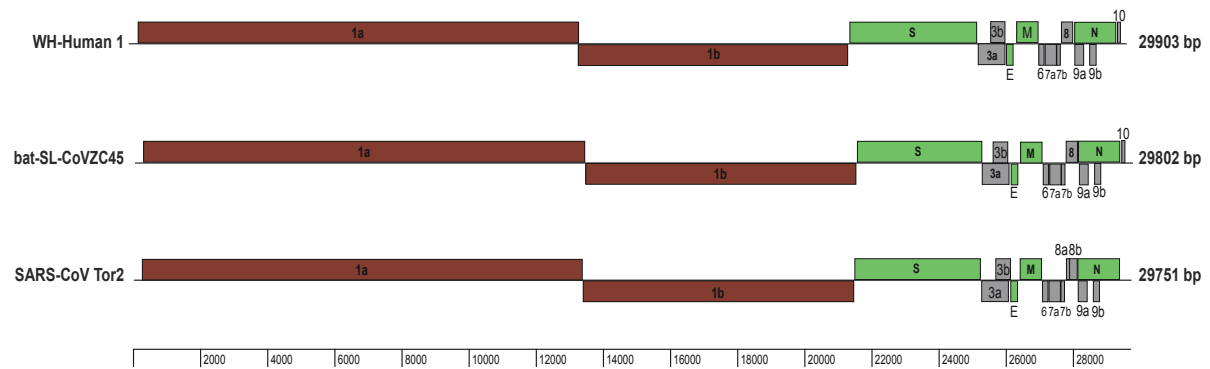
Characteristic	Patient 1	Patient 2	Patient 3	Patient 4	Patient 5	Patient 6	Patient 7
Age (Year)	41	44	42	70	31	51	43
Sex	M	F	M	F	M	M	M
Date of illness onset	Dec 20,2019	Dec 22,2019	Dec 24,2019	Dec 24,2019	Dec 21,2019	Dec 16,2109	Dec 14,2019
Date of admission	Dec 26,2019	Dec 22,2019	Dec 28,2019	Dec 28,2019	Dec 28,2019	Dec 27,2019	Dec 14,2019
Fever	+	+	+	+	+	+	+
Body Temperature (°C)	38.4	37.3	39	37.9	38.7	37.2	38
Cough	+	+	+	+	+	+	+
Sputum Production	+	+	-	-	-	-	+
Dizzy	+	-	+	+	-	+	-
Weakness	+	-	+	-	-	-	-
Chest tightness	+	-	-	+	+	-	-
Dyspnea	+	-	-	+	+	+	-
Bacterial culture	-	-	streptococcus pneumoniae	streptococcus pneumoniae	-	-	streptococcus pneumoniae
Glucocorticoid therapy	No	No	Yes	Yes	Yes	Yes	No
Antibiotic therapy	Cefoselis	Ceftazidime, Levofloxacin	Cefminox	Cefminox, moxifloxacin	Cefminox	No	No
Antiviral therapy	Oseltamivir	No	Oseltamivir, ganciclovir	Oseltamivir, ganciclovir	Oseltamivir, ganciclovir	Oseltamivir, ganciclovir	No
Oxygen therapy	mechanical ventilation	No	No	Mask	No	No	No



480

481 **Figure 1.** Chest radiographs of patient 1. **a. b. c. d.** Chest computed tomographic scans of
482 Patient 1 were obtained on the day of admission (day 6 after the onset of illness). Bilateral
483 focal consolidation, lobar consolidation, and patchy consolidation were clearly observed,
484 especially in the lower lung. **e.** Chest radiograph of patient 1 was obtained on day 5 after
485 admission (day 11 after the onset of illness). Bilateral diffuse patchy and fuzzy shadow were
486 observed.

487

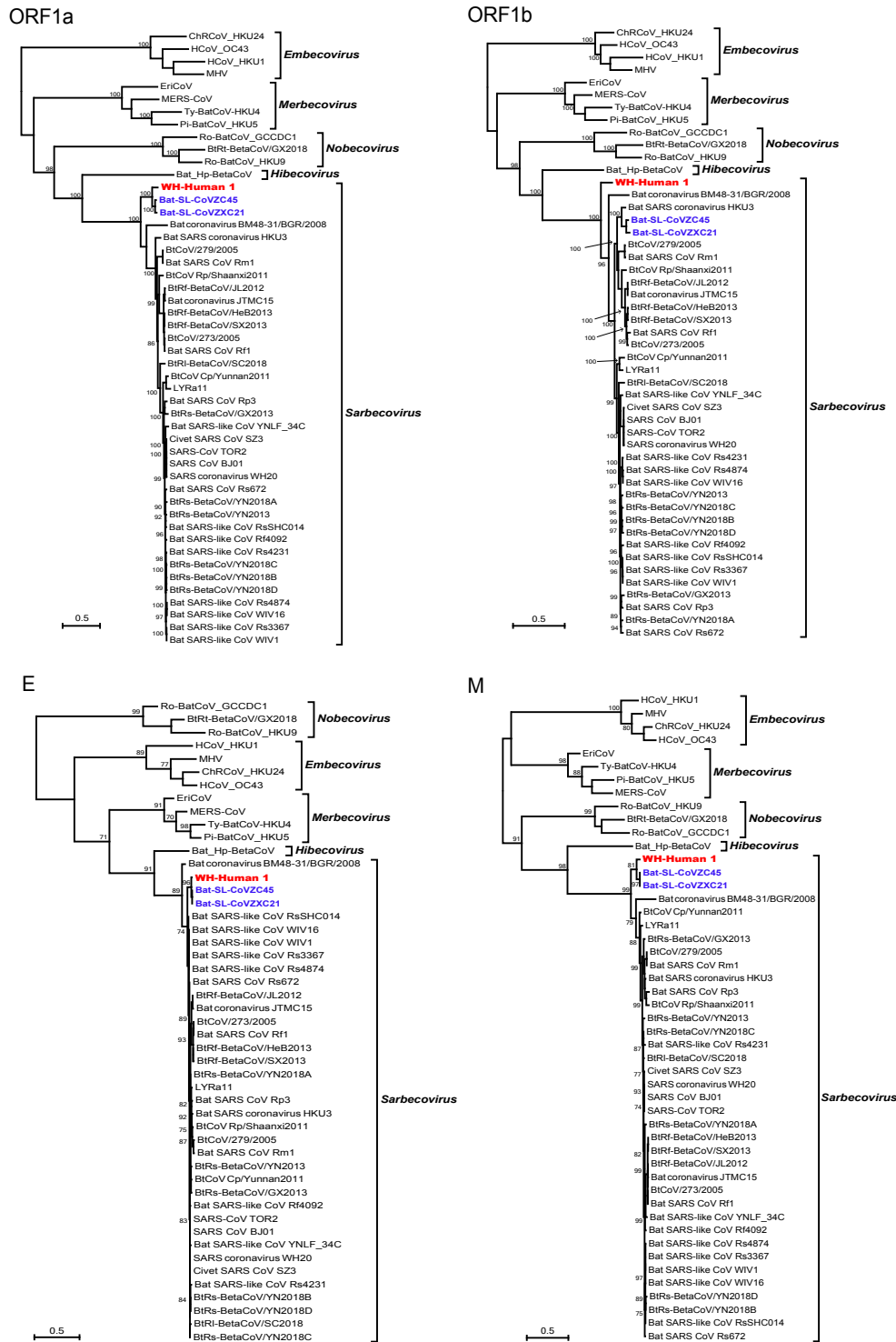


488

489 **Figure 2.** Genome organization of SARS and SARS-like CoVs including Tor2, CoVZC45

490 and WHCV determined here.

491

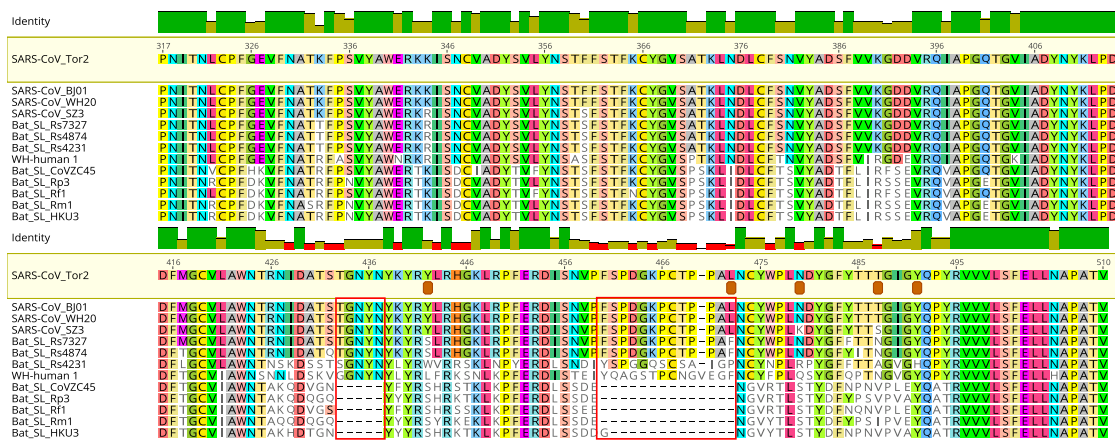


492

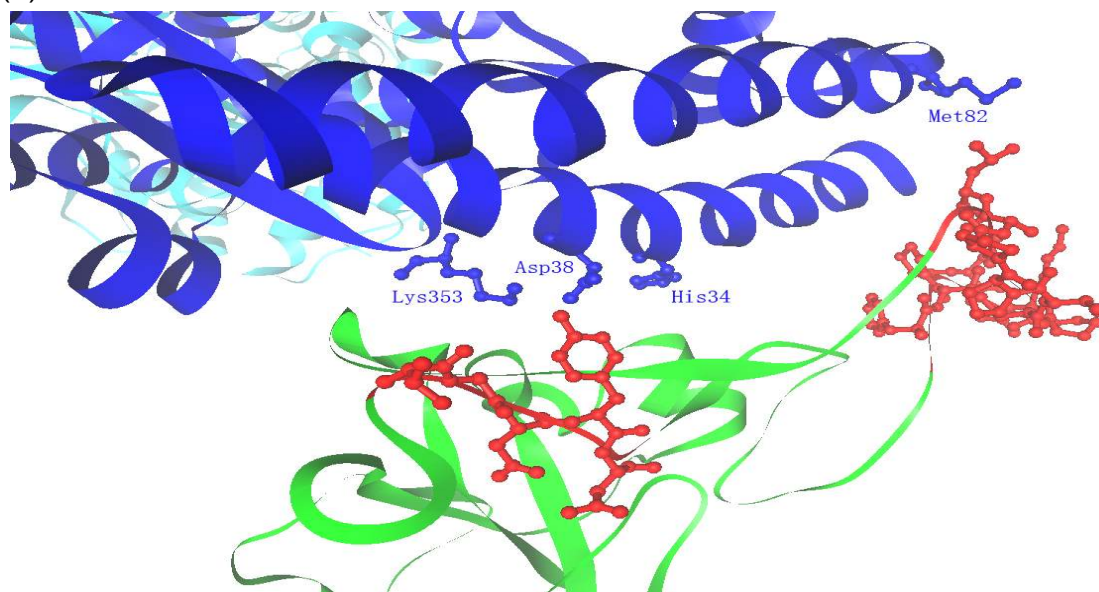
493 **Figure 3.** Maximum likelihood phylogenetic trees of nucleotide sequences of the ORF1a,
 494 ORF1b, E and M genes of WHCV and related coronaviruses. Numbers (>70) above or below
 495 branches indicate percentage bootstrap values for the associated nodes. The trees were mid-
 496 point rooted for clarity only. The scale bar represents the number of substitutions per site.

497

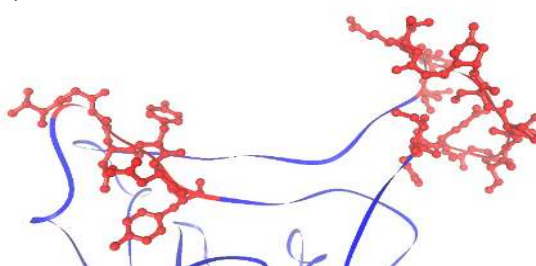
(a)



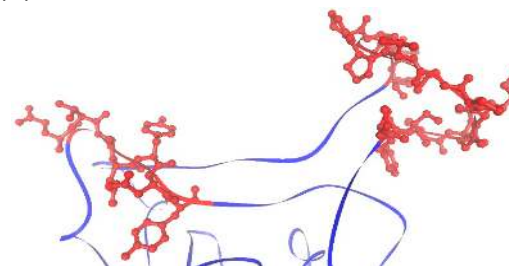
(b)



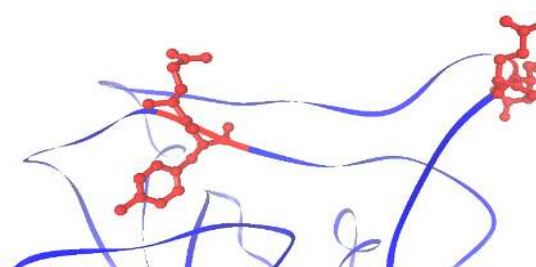
(c)



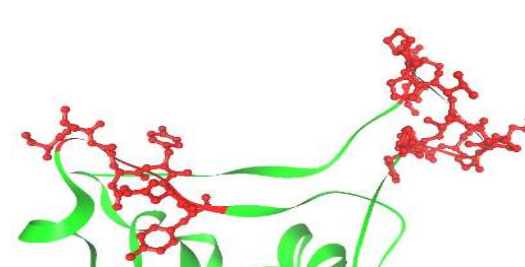
(d)



(e)



(f)



499 **Figure 4.** Analysis of receptor-binding domain (RBD) of the spike (S) protein of WHCV
500 coronavirus. **(a)** Amino acid sequence alignment of SARS-like CoV RBD sequences. Three
501 bat SARS-like CoVs, which could efficiently utilize the human ACE2 as receptor, had an
502 RBD sequence of similar size to SARS-CoV, and WHCV contains a single Val 470 insertion.
503 The key amino acid residues involved in the interaction with human ACE2 are marked with a
504 brown box. In contrast, five bat SARS-like CoVs had amino acid deletions at two motifs
505 (amino acids 473-477 and 460-472) compared with those of SARS-CoV, and Rp3 has been
506 reported not to use ACE2.¹¹ **(b)** The two motifs (aa 473-477 and aa 460-472) are shown in red
507 on the crystal structure of the SARS-CoV spike RBD complexed with receptor human ACE2
508 (PDB 2AJF). Human ACE2 is shown in blue and the SARS-CoV spike RBD is shown in
509 green. Important residues in human ACE2 that interact with SARS-CoV spike RBD are
510 marked. **(c)** Predicted protein structures of RBD of WHCV spike protein based on target-
511 template alignment using ProMod3 on the SWISS-MODEL server. The most reliable models
512 were selected based on GMQE and QMEAN Scores. Template: 2ghw.1.A, GMQE: 0.83;
513 QMEAN:-2.67. Motifs resembling amino acids 473-477 and 460-472 of the SARS-CoV spike
514 protein are shown in red. **(d)** Predicted structure of RBD of SARS-like CoV Rs4874.
515 Template: 2ghw.1.A, GMQE:0.99; QMEAN:-0.72. Motifs resembling amino acids 473-477
516 and 460-472 of the SARS-CoV spike protein are shown in red. **(e)** Predicted structure of the
517 RBD of SARS-like CoV Rp3. Template: 2ghw.1.A, GMQE:0.81, QMEAN:-1.50. **(f)** Crystal
518 structure of RBD of SARS-CoV spike protein (green) (PDB 2GHV). Motifs of amino acids
519 473-477 and 460-472 are shown in red.
520

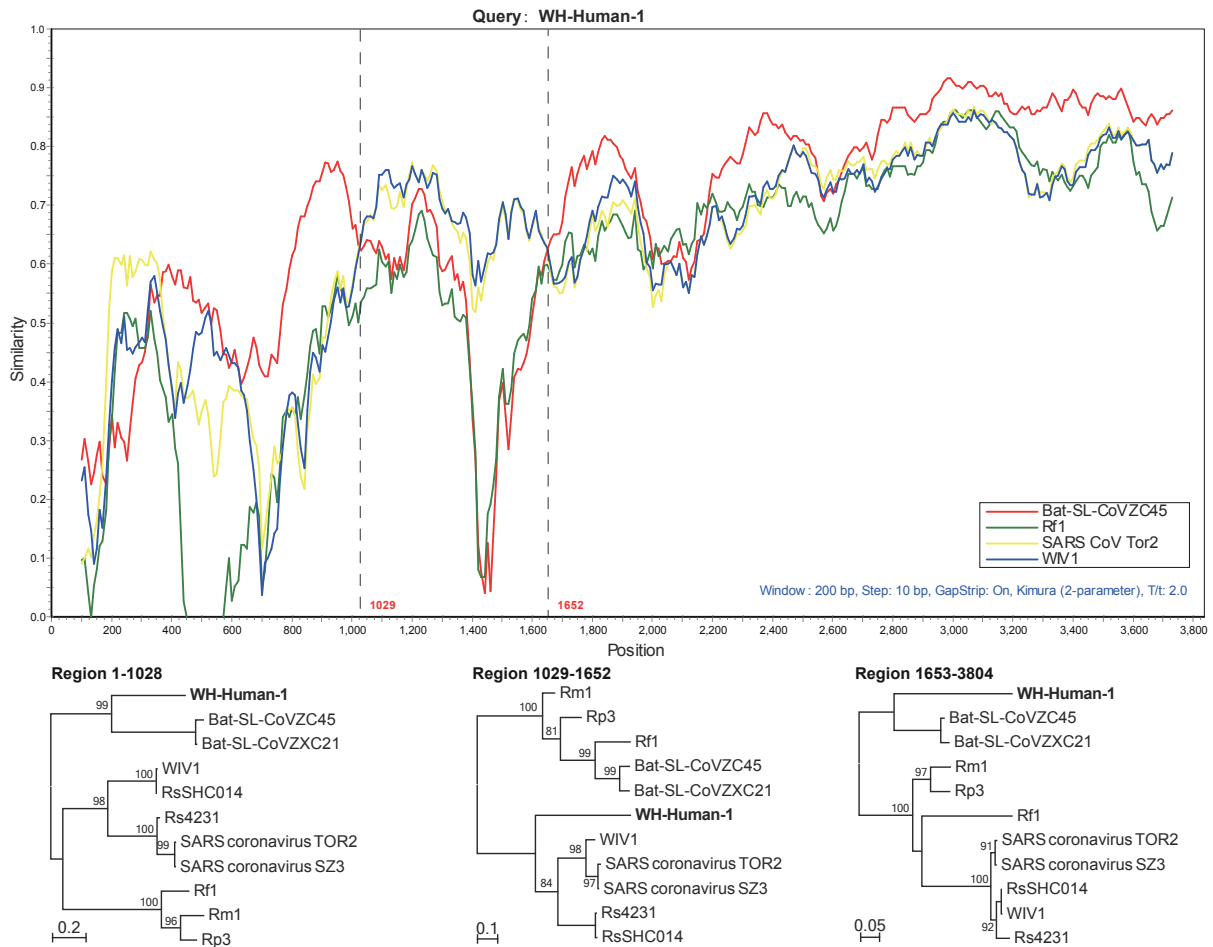


Figure 5. Possible recombination events in the S gene of sarbecoviruses. A sequence similarity plot (upper panel) reveals two putative recombination break-points shown by black dashed lines, with their locations indicated at the bottom. The plot shows S gene similarity comparisons of the WHCV (query) against SARS-CoV Tor2 and bat SARS-like CoVs WIV1, Rf1 and CoVZC45. Phylogenies of major parental region (1-1028 and 1653-3804) and minor parental region (1029-1652) are shown below the similarity plot. Phylogenies were estimated using a ML method and mid-point rooted for clarity only. Numbers above or below branches indicate percentage bootstrap values.