

1

2 **The codon usage pattern of the novel coronavirus is**

3 **drastically different from those of other pathogenic viruses**

4 *Xiaolong Wang**

5 *College of Life Sciences, Ocean University of China, Qingdao, 266003, P. R. China*

6 **Abstract**

7 The current outbreak of a novel coronavirus (COVID-19) has caused thousands of deaths and
8 has been declared to be a worldwide pandemic by the World Health Organization. There have been
9 various disputes but the origin of COVID-19 is not clear. Here we analyzed the similarities of codon
10 usage patterns between humans and pathogenic viruses, such as human immunodeficiency virus
11 (HIV), highly pathogenic avian influenza (HPAI), SARS, MERS, and COVID-19. In HIVs, HPAIs,
12 SARS, and MERS, codon usages are highly similar to that of humans; in contrast, the codon usage
13 pattern of COVID-19 is drastically different from those of humans and other pathogenic viruses.
14 Besides, coronaviruses have been evolving in two opposite directions: human-preferred codons are
15 adopted to substitute less-preferred ones in SARS and MERS but are substituted by less-preferred
16 ones in COVID-19. The unique codon usage pattern suggesting that COVID-19 was evolved in an
17 intermediate host, in which its codon usage pattern becomes drastically different from that of bats
18 or humans, and its pathogenicity is weakened compared with SARS and MERS COVs. Finally, we
19 appeal to international cooperation to eliminate the epidemic by cutting off the transmission routes
20 among humans and to search for the origin and intermediate hosts of the novel coronavirus to
21 prevent future animal-to-human transmission.

22

23

24

25

26

27

28

29 *Xiaolong Wang, College of Life Sciences, Ocean University of China, No. 5 Yushan Road,
30 Qingdao, 266003, Shandong, China, E-mail: Xiaolong@ouc.edu.cn.

1

2 **1. Background**

3 In the last two decades, three serious epidemics caused by pathogenic coronavirus have
4 emerged, including Severe Acute Respiratory Syndrome (SARS) in 2002-2003 [1], the Middle East
5 Respiratory Syndrome (MERS) in 2012-2015 [2], and the current outbreak of a novel coronavirus
6 (COVID-19). COVID-19 has caused thousands of deaths and hundreds of thousands of hospitalized
7 cases not only in China but at present more seriously in all over the world and has been declared to
8 be a worldwide pandemic by the World Health Organization.

9 Most pathogenic viruses are of zoonotic origin. For example, human immunodeficiency virus
10 (HIV) was originated from the chimpanzee simian immune deficiency virus (SIVcpz) [3], highly
11 pathogenic avian influenza (HPAI) was originated from bird influenza [4], and pathogenic COVs
12 are originated from a bat coronavirus [5-7]. Lentivirus like HIV inhabits in a host with no symptom
13 for a long period; coronaviruses, such as SARS and MERS, cause severe acute immune responses
14 and respiratory infections in a short period, may cause the death of the host if the viruses are not
15 eliminated by the immune system or medical treatment.

16 Viral genomes are small in size and largely rely on the host to execute biological activities like
17 replication, protein synthesis, and transmission. After the invasion of a human body, viruses adjust
18 their growth rates and change their pathogenicity/immunogenicity to adapt for a short- or long-term
19 inhabiting in humans. A common strategy for the evolution of viruses is to change the usage of
20 codons, which has strong impacts on viral gene expression and the progress of the pathogenic virus.
21 In 1996, Haas, Park, and Seed reported that the change of codon usage can lead to the inhibition of
22 HIV protein synthesis and the limitation in the expression of HIV-1 envelop glycoprotein [8]. In
23 2017, Roy, Banerjee, and Basak demonstrated that the rate of substitution in the envelop gene is
24 associated with disease progression [9]. Also, it was suggested that mutational pressure, rather than
25 natural selection for specific coding triplets, is the main determinant of codon usage [10]. In 2013,
26 Moratorio and his colleagues did a comprehensive analysis of the West Nile virus (WNV), which
27 suggested that the genomic biases are the result of the evolution of genome composition, the need
28 to escape the antiviral cell responses and to re-adapt its codon usage to different environments [11].

29 To analyze the evolutionary characteristics of the novel coronaviruses, we analyzed the codon

1 usages of SARS, MERS, and COVID-19, determined the changes of the codon usages by compared
2 with those of their most recent common ancestors and those of other pathogenic viruses, including
3 HIVs and HPAs. The codon usages of the coronaviruses are associated with their high growth rates
4 and severe acute inflammatory responses in humans, provides a theoretical basis for the prediction
5 of possible changes of the pathogenic coronaviruses in the future.

6 **2. Methods**

7 **2.1 Genomes sequences**

8 The reference genome sequences and all available complete genome sequences of COVID-19,
9 SARS, MERS, HIV, SIVcpz, and HPAI were downloaded from the NCBI Nucleotide Database
10 during March 1st-16th, 2020. The accession numbers of the reference genome sequences are:
11 NC_004718.3 (SARS), NC_019843.3 (MERS), NC_045512.2 (COVID-19), NC_001802.1 (HIV1),
12 AF115393.1 (SIVcpz), NC_002022.1 (H1N1), NC_007361.1 (H5N1), NC_026422.1 (H7N9) and
13 AF250131.1 (H7N2), respectively.

14 **2.2 Phylogenetic trees**

15 We constructed a multiple sequence alignment of 299 complete coronavirus genomes of using
16 a phylogeny-aware alignment software, PRANK v170427. Maximum likelihood phylogenies were
17 estimated using PhyML v3.115, utilizing the GTR+I+G model of nucleotide substitution with 1,000
18 bootstrap replicates. The phylogenetic tree was plotted using MEGA v7.0.26 [12].

19 **2.3 Analyze of codon usages**

20 Genes and genomes display a non-random usage of synonymous codons for specific amino
21 acids. A measure of the extent of this non-randomness is given by the relative synonymous codon
22 usage (RSCU), which is calculated as the ratio of the observed frequency of the codons divided by
23 the expected frequency of the same codon if codon usage was uniform within a synonymous codon
24 group [13]. An RSCU value greater than one indicates that the observed frequency of synonymous
25 codons is more preferred compared to the expected frequency [14]. RSCU values of the 59 codons
26 [excluding the single synonymous codons, AUG (Met) and UGG (Trp) and the termination codons,
27 UGA, UAG, and UAA] of all coding gene sequences were calculated using CodonW v 1.4.2.

28 **2.4 Assessment of the distance and the similarity index of codon usages**

29 The relationship among the codon usages of humans and different viruses was calculated using

1 a squared Euclidean distance method as described by Wei Ji *et al* [15], which is computed as follows:

2

3

$$d(H, V) = \sum_{i=1}^{59} (h_i - v_i)^2$$

4 where $d(H, V)$ represents the distance between the overall codon usage pattern of human and
 5 a specific virus, h_i indicates the RSCU value for a particular codon in human, v_i signifies the RSCU
 6 value of the same codon for a certain viral gene or genome.

7 We also used a similarity index of the codon usages, as described by Roy, Banerjee, and Basak
 8 [9], to understand the influence of the host genome on the adaptability of the virus genome inside
 9 the host. The influence of the overall codon usage pattern of the host on the formation of the codon
 10 usage of the virus is defined as the similarity index, which is computed as follows:

$$11 \quad R(H, V) = \frac{\sum_{i=1}^{59} h_i \cdot v_i}{\sqrt{\sum_{i=1}^{59} h_i^2 \cdot \sum_{i=1}^{59} v_i^2}}$$

$$12 \quad D(H, V) = \frac{1 - R(H, V)}{2}$$

13 where $R(H, V)$ represents the degree of similarity between the overall codon usage pattern of
 14 human (H) and that of a specific viral gene/genome (V), h_i indicates the RSCU value for a particular
 15 codon in human, v_i signifies the RSCU value of the same codon for a certain viral gene/genome. D
 16 (H, V) represents the potential effect of the overall codon usage of humans on that of the virus. This
 17 value ranges from 0.0 to 1.0 and useful for cross-species comparison of codon usages.

18 **2.5 Codon and aa unified sequence alignment**

19 The packaging and the fusion of a virus into a cell rely on their surface/envelop proteins. The
 20 spike glycoprotein (S) of COVs and the envelop glycoproteins (GP120) of HIV have become the
 21 first choice of the targets in various studies. Here, the *gp120* gene of HIVs and the *s* gene of COVs
 22 were aligned by Codon-AA Unified Sequence Alignment (CAUSA v2.1.018) [16]. By comparing
 23 with their most recent common ancestors, synonymous and nonsynonymous codon substitutions
 24 were found and subject to the analyses of the change of codon preferences.

25 **3. Results**

1 **3.1 Phylogenetic analyses**

2 All available complete genome sequences that are related to SARS and COVID-19 viruses
3 were aligned and a genome-wide maximum likelihood phylogenetic tree was established by phyML.
4 As of [March 1st, 2020](#), there were [45](#) COVID-19 viral genomes deposited in GenBank. As shown
5 in [Fig 1a](#), the overall phylogenetic tree is consistent with those reported earlier [15, 17-19]. COVID-
6 19 share 79.5% identify to SARS-Cov and is 96% identical to a bat coronavirus (RaTG13) at the
7 whole genome level [18], which is identified as the most recent common ancestor of COVID-19
8 and SARS COVs.

9 **3.2 Codon usages of different viruses and their similarities to that of humans**

10 The RSCU values of different viral genomes were compared with that of humans to assess the
11 influence of the human host in shaping the patterns of codon usage among the viruses. It has been
12 reported that the rate of codon substitution in the envelop gene is associate with disease progression,
13 differs among the three different types of HIV, rapid progressor (RP), slow progressor (SP), and
14 long-term non-progressor (LTNP) of HIV1 infected individuals [9]. Based on the RSCU values for
15 different viruses given by CodonW, the relationship among codon usages of humans and different
16 viruses was calculated using a squared Euclidean distance and a similarity index of codon usages.
17 As shown in [Table 1](#), the codon usage patterns of HIVs are all similar to that of human and that of
18 HPAs are even more similar to that of humans. The codon usage patterns of SARS and MERS are
19 also highly similar to that of humans, however, COVID-19 has a very special codon usage pattern
20 which is drastically different from that of humans, suggesting that COVID-19 was evolved in an
21 intermediate host, in which its codon usage pattern becomes drastically different from that of bats
22 or humans, and its pathogenicity is significantly weakened compared with SARS ad MERS COVs.
23 Recently, it is reported that the intermediate hosts could be snakes [15] or pangolins [17], but further
24 investigations are needed to validate these speculations.

25 **3.3 The changes of codon usages in the protein-coding genes**

26 Because the sizes of viral genomes are very small, the differences of codon usages could be
27 obscured by noise when they were calculated by counting the number of codons used in the genome
28 sequences. As shown in [Fig 2](#), we performed codon alignments of their surface/envelop proteins,
29 identified synonymous and nonsynonymous codon substitutions, calculate the codon preferences,

1 and investigated whether codon preferences have been changed in different viruses.

2 When compared with the *gp120* gene of SIVcpz, the HIV *gp120* gene has 226 synonymous
3 and 279 nonsynonymous substitutions. The average RSCU of the substitutional codons is used as
4 an index for human preference (HPI). As shown in Table 1, HPI decreased in the nonsynonymous
5 substitutions significantly (paired t-test $P=0.0301$). In contrast, compared with the *s* gene of the bat
6 coronavirus RaTG13, the *s* gene of SARS contains 450 synonymous and 262 nonsynonymous
7 substitutions, while that of COVID-19 contains only 215 synonymous and 29 nonsynonymous
8 substitutions. As shown in Table 1, compared with RaTG13, HPI increased in SARS but decreased
9 in COVID-19.

10 Besides, compared with SARS, HPI decreased even further in COVID-19 in both synonymous
11 and nonsynonymous substitutions. Although neither of the differences of preference among SARS,
12 COVID-19 and bat COV is statistically significant, the difference of the preference of the
13 synonymous substitutions between COVID-19 and SARS is close to statistically significant (paired
14 t-test $P=0.0523$). It is clear that coronaviruses are evolving in two opposite directions: in SARS,
15 human-preferred codons are adopted to substitute less-preferred ones; in COVID-19, however,
16 human-preferred codons are abandoned and substituted by less-preferred ones, suggesting that the
17 HPI of COVID-19 has been decreasing since it was isolated from bat COV.

18 2. Discussion & Conclusion

19 The above analysis concludes that codon usages have been changed in tested human pathogenic
20 viruses comparing with their ancestors in wild animals. In HIV, HPAs, SARS, and MERS, codon
21 usages are highly similar to that of humans. In contrast, in COVID-19, hundreds of human-preferred
22 codons were substituted by synonymous codons that are less preferred in humans, making its codon
23 usage patterns drastically different from that of humans.

24 Moreover, SARS and MERS have an excessive number of highly human-preferred codons, the
25 growth rate of them will be too fast and dysregulated in an infected human body, rob host cells of
26 too many nutrients, energy, and resources. After infection, the fast growth of viruses is the cause
27 of high mortality of the patients, as it may trigger a severe acute response, an inflammatory storm in
28 the human body. Compared with the codon usages of SARS/MERS, the codon usage of COVID-
29 19 is more different from that of humans. On one hand, COVID-19 infection is therefore not as

1 severe as SARS and MERS, on the other hand, however, as it is much milder than SARS and MERS,
2 COVID-19 is indeed a more successful pathogenic coronavirus and perhaps has greater potential.

3 Like other pathogenic viruses, the coronaviruses evolve by optimizing the sequence, structure,
4 and functionality of their proteins by changing their codons. If the current epidemic could not be
5 eliminated in a short period, very likely, the coronavirus will develop a chronic disease eventually.
6 It may evolve either into an HIV-like lentivirus or a flu-like self-limiting virus, or both, but they
7 may keep their severe acute pathogenicity and remain to be dangerous for a long period. As a novel
8 pathogenic coronavirus, they are in the early stage of their evolutionary journey in humans. . Finally,
9 we appeal to international cooperation to eliminate the epidemic by cutting off the transmission
10 routes among humans and to search for the origin and intermediate hosts of the novel coronavirus
11 to prevent future animal-to-human transmission.

12 **Data availability**

13 This study conduct data analyses based on existing gene/genome sequences that are available
14 in the NCBI Nucleotide Database and the Global Initiative on Sharing Avian Influenza Data
15 (GISAID) database. The list of GenBank accession numbers of the genome sequences is available
16 online as a text file (AllCoronaVirus.list.txt). The RSCU data for human and viruses are available
17 online as a excel spreadsheet (RSCU-human-viruses.xlsx).

18 **Acknowledgments**

19 This research is funded by the National Natural Science Foundation of China (Grant 31571369).
20 We acknowledge the doctors, nurses, and scientists from all over the world for battling against
21 the epidemic and making the genomic sequences of coronavirus freely and publicly available.
22

23 **Competing interests**

24 We declare that we have no conflicts of interest.

25 **References**

26

- 27 1. Lau, S.K., et al., *Severe acute respiratory syndrome coronavirus-like virus in Chinese*
28 *horseshoe bats*. Proc Natl Acad Sci U S A, 2005. **102**(39): p. 14040-5.
- 29 2. Zaki, A.M., et al., *Isolation of a novel coronavirus from a man with pneumonia in Saudi*
30 *Arabia*. N Engl J Med, 2012. **367**(19): p. 1814-20.

- 1 3. Rambaut, A., et al., *Human immunodeficiency virus. Phylogeny and the origin of HIV-1.*
2 Nature, 2001. **410**(6832): p. 1047-8.
- 3 4. Lee, D.H., et al., *Evolution, global spread, and pathogenicity of highly pathogenic avian*
4 *influenza H5Nx clade 2.3.4.4.* J Vet Sci, 2017. **18**(S1): p. 269-280.
- 5 5. Hon, C.C., et al., *Evidence of the recombinant origin of a bat severe acute respiratory*
6 *syndrome (SARS)-like coronavirus and its implications on the direct ancestor of SARS coronavirus.* J
7 Virol, 2008. **82**(4): p. 1819-26.
- 8 6. Mohd, H.A., J.A. Al-Tawfiq, and Z.A. Memish, *Middle East Respiratory Syndrome*
9 *Coronavirus (MERS-CoV) origin and animal reservoir.* Virol J, 2016. **13**: p. 87.
- 10 7. Hu, B., et al., *Discovery of a rich gene pool of bat SARS-related coronaviruses provides new*
11 *insights into the origin of SARS coronavirus.* PLoS Pathog, 2017. **13**(11): p. e1006698.
- 12 8. Haas, J., E.C. Park, and B. Seed, *Codon usage limitation in the expression of HIV-1 envelope*
13 *glycoprotein.* Curr Biol, 1996. **6**(3): p. 315-24.
- 14 9. Roy, A., R. Banerjee, and S. Basak, *HIV Progression Depends on Codon and Amino Acid*
15 *Usage Profile of Envelope Protein and Associated Host-Genetic Influence.* Front Microbiol, 2017. **8**:
16 p. 1083.
- 17 10. Shackelton, L.A., C.R. Parrish, and E.C. Holmes, *Evolutionary basis of codon usage and*
18 *nucleotide composition bias in vertebrate DNA viruses.* J Mol Evol, 2006. **62**(5): p. 551-63.
- 19 11. Moratorio, G., et al., *A detailed comparative analysis of the overall codon usage patterns*
20 *in the West Nile virus.* Infect Genet Evol, 2013. **14**: p. 396-400.
- 21 12. Kumar, S., G. Stecher, and K. Tamura, *MEGA7: Molecular Evolutionary Genetics Analysis*
22 *Version 7.0 for Bigger Datasets.* Mol Biol Evol, 2016. **33**(7): p. 1870-4.
- 23 13. Sharp, P.M. and W.H. Li, *An evolutionary perspective on synonymous codon usage in*
24 *unicellular organisms.* J Mol Evol, 1986. **24**(1-2): p. 28-38.
- 25 14. dos Reis, M., L. Wernisch, and R. Savva, *Unexpected correlations between gene expression*
26 *and codon usage bias from microarray data for the whole Escherichia coli K-12 genome.* Nucleic
27 Acids Res, 2003. **31**(23): p. 6976-85.
- 28 15. Ji, W., et al., *Cross-species transmission of the newly identified coronavirus 2019-nCoV.* J
29 Med Virol, 2020. **92**(4): p. 433-440.
- 30 16. Wang, X. and C. Yang, *CAUSA 2.0: accurate and consistent evolutionary analysis of proteins*
31 *using codon and amino acid unified sequence alignments.* PeerJ PrePrints **3**.
- 32 17. Lam, T.T.-Y., et al., *Identification of 2019-nCoV related coronaviruses in Malayan pangolins*
33 *in southern China.* bioRxiv, 2020: p. 2020.02.13.945485.
- 34 18. Zhou, P., et al., *A pneumonia outbreak associated with a new coronavirus of probable bat*
35 *origin.* Nature, 2020. **579**(7798): p. 270-273.
- 36 19. Matsuda, T., H. Suzuki, and N. Ogata, *Phylogenetic analyses of the severe acute respiratory*
37 *syndrome coronavirus 2 reflected the several routes of introduction to Taiwan, the United States,*
38 *and Japan.* 2020: arXiv:2002.08802 [q-bio.GN].
- 39
- 40

1 Table 1. Codon usages of different type of virus and their distance/similarity to that of humans
 2

Virus	Strain / Type	<i>d</i> (H, V)	<i>R</i> (H, V)	<i>D</i> (H, V)
SIV	SIVcpz	27.1721	0.8263	0.0868
	RP	27.6765	0.8176	0.0912
HIVs	SP	28.2044	0.8143	0.0929
	LTNP	28.3472	0.8145	0.0927
HPAIs	H5N1	12.8378	0.9116	0.0442
	H7N9	13.0191	0.9069	0.0466
	H1N1	13.0251	0.9066	0.0467
	H7N2	14.4611	0.8972	0.0514
COVs	MERS	26.1475	0.8301	0.0850
	SARS	27.9605	0.8245	0.0878
	COVID-19	39.3928	0.7636	0.1182

3
 4

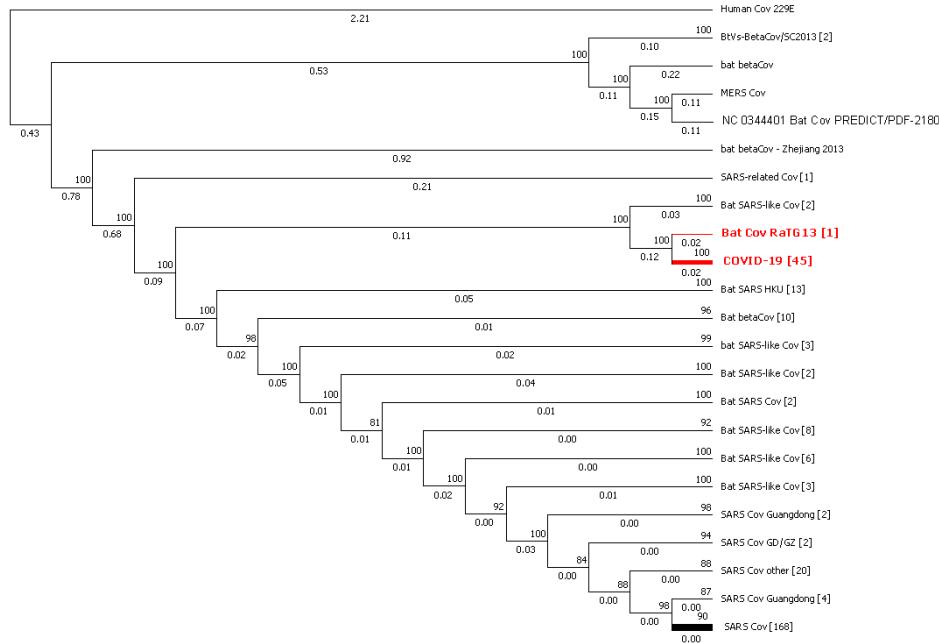
1 Table 2. Codon substitutions and the average RSCU frequencies of codons of the envelop or surface protein
 2

Compare	Type of codon Substitution	Number of Codon Substitutions	Human Preference Index (HPI) (Average RSCU Frequency Per Thousand)		
			SIVcpz	HV1	P-value
<i>HIV1 vs SIV</i>	Synonymous	226	17.7783	18.2407 ↑	0.2877
	Nonsynonymos	279	18.4039	17.1276 ↓	0.0301*
			Bat Cov	SARS	P-value
<i>SARS vs Bat Cov</i>	Synonymous	450	16.7153	17.0744 ↑	0.1856
	Nonsynonymos	262	16.8462	17.2905 ↑	0.2242
			Bat Cov	Cov-2019	P-value
<i>Covid-19 vs Bat Cov</i>	Synonymous	215	17.7167	17.2065 ↓	0.2052
	Nonsynonymos	29	16.9897	15.6690 ↓	0.2200
			SARS	COVID-9	P-value
<i>Covid-19 vs SARS</i>	Synonymous	435	17.2218	16.5579 ↓	0.0523
	Nonsynonymos	265	17.2400	16.6743 ↓	0.1710

3
 4

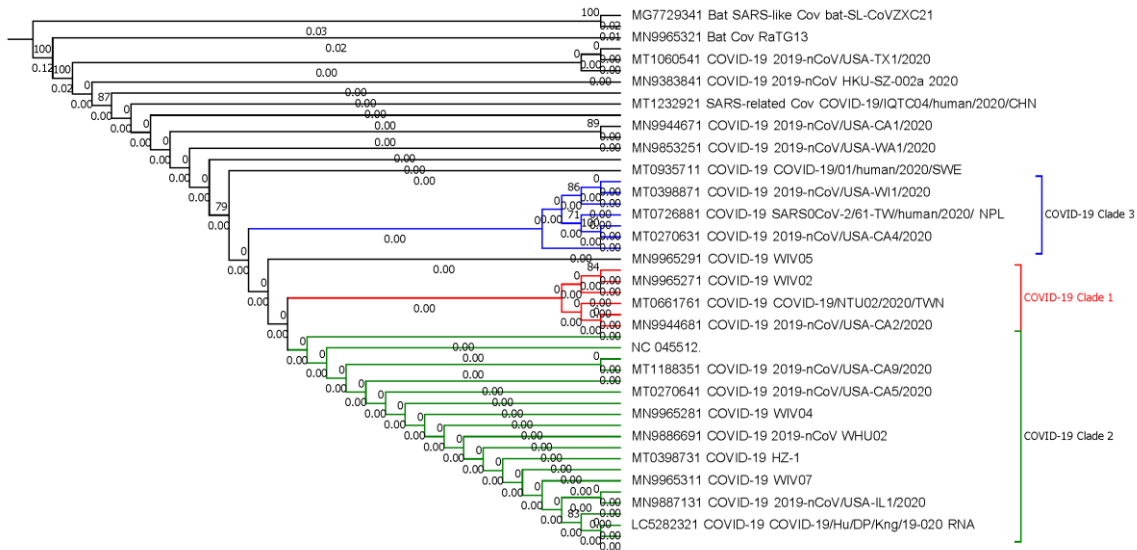
1
2

A



3
4

B



5 **Fig 1. The molecular phylogenetic tree of coronavirus.** (A) The phylogenetic tree of all SARS-related coronavirus; (B)
 6 the subtree of the novel coronavirus (COVID-19) in A. The analysis base on a PRANK alignment of 299 complete
 7 coronavirus genome sequences. All positions containing gaps and missing data were eliminated. The evolutionary history
 8 was inferred by using the Maximum Likelihood method by phyML. The tree with the highest log is shown. Trees were
 9 plotted in MEGA7. Root was placed on the most distant branch, human COV 229E. The bootstrap percentage of trees in
 10 which the associated taxa clustered together is shown next to the branches. The distances between the taxa are shown in
 11 the middle of the branches.

1
2
3

A

Seq	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65
HV1J3	Ctgt	Sagl	Agct	Agca	Egaa	Qcaa	Lttg	Wggg	Vgtc	Taca	Vgtc	Ytat	Ytat	Gggg	Vgta	Pcct	Vgtg	Wggg	Kaaa	Egat	Agca	Agcc	Tacc	Tact	Lcta	Fttt	Ctgt	Agca	Stca	Dgat	Agct
HV1ZH	Ctgt	Sagl	Agct	Agca	Egaa	Naac	Lttg	Wggg	Vgtc	Tacc	Vgtc	Ytac	Ytat	Gggg	Vgta	Pcct	Vgtg	Wggg	Kaaa	Dgat	Agca	Egat	Tacc	Tact	Lcta	Fttt	Ctgt	Agca	Stca	Dgat	Agct
HV1B1	Ctgt	Sagl	Agct	Taca	Egaa	Kaaa	Lttg	Wggg	Vgtc	Taca	Vgtc	Ytat	Ytat	Gggg	Vgta	Pcct	Vgtg	Wggg	Kaag	Egat	Agca	Tacc	Tacc	Tact	Lcta	Fttt	Ctgt	Agca	Stca	Dgat	Agct
HV1B1	Ctgt	Sagl	Agct	Taca	Egaa	Kaaa	Lttg	Wggg	Vgtc	Taca	Vgtc	Ytat	Ytat	Gggg	Vgta	Pcct	Vgtg	Wggg	Kaag	Egat	Agca	Tacc	Tacc	Tact	Lcta	Fttt	Ctgt	Agca	Stca	Dgat	Agct
HV1A2	Ctgt	Sagl	Agct	Taca	Egaa	Kaaa	Lttg	Wggg	Vgtc	Taca	Vgtt	Ytat	Ytat	Gggg	Vgta	Pcct	Vgtg	Wggg	Kaaa	Egat	Agca	Tact	Tacc	Tact	Lcta	Fttt	Ctgt	Agca	Stca	Dgat	Agct
HV10Y	Ctgt	Sagl	Agct	Agca	Egaa	Naat	Lttg	Wggg	Vgtc	Taca	Vgtc	Ytat	Ytat	Gggg	Vgta	Pcct	Vgtg	Wggg	Kaaa	Egat	Agca	Tacc	Tacc	Tact	Lcta	Fttt	Ctgt	Agca	Stca	Dgat	Agct
HV1RH	Ctgt	Sagl	Agct	Agca	Egat	Dgac	Lttg	Wggg	Vgtc	Taca	Vgtc	Ytat	Ytat	Gggg	Vgta	Pcct	Vgtg	Wggg	Kaaa	Egat	Agca	Tacc	Tacc	Tact	Lcta	Fttt	Ctgt	Agca	Stca	Egat	Agct
HV1C4	Ctgt	Sagl	Agct	Agca	Agca	Naac	Lttg	Wggg	Vgtc	Taca	Vgtc	Ytat	Ytat	Gggg	Vgta	Pcct	Vgtg	Wggg	Kaaa	Egat	Agca	Tacc	Tacc	Tact	Lcta	Fttt	Ctgt	Agca	Stca	Dgat	Agct
HV1EL	Ctgt	Sagl	Agct	Agca	Dgac	Naat	Lctg	Wggg	Vgtc	Taca	Vgtt	Ytat	Ytat	Gggg	Vgtg	Pcct	Vgta	Wggg	Kaag	Egat	Agca	Tacc	Tacc	Tact	Lcta	Fttt	Ctgt	Agca	Stca	Dgat	Agct
HV1ND	Ctgt	Sagl	Agct	Agca	Egat	Dgat	Lttg	Wggg	Vgtc	Taca	Vgtt	Ytat	Ytat	Gggg	Vgtg	Pcct	lata	Wggg	Kaag	Egat	Agca	Tact	Tacc	Tact	Lcta	Fttt	Ctgt	Agca	Stca	Dgat	Agct
HV1MA	Ctgt	Sagl	latt	Agca	Egat	Dgat	Lttg	Wggg	Vgtt	Taca	Vgtt	Ytat	Ytat	Gggg	Vgta	Pcct	Vgtg	Wggg	Kaaa	Egat	Agca	Tacc	Tact	Tact	Lcta	Fttt	Ctgt	Agca	Stca	Dgat	Agct
SIVCZ	Ctgt	Lttg	Tacc	Stcl	Eg--	--ag	Ltta	Wggg	Vgta	Taca	Vgta	Ytat	Ytat	Ggga	Vgta	Pcct	Vgtt	Wggg	Hcat	Dgat	Agct	Dgac	Pccg	Vgta	Lctc	Fttt	Ctgt	Agcc	Stca	Dgac	Agct

4
5
6

B

Seq	670	671	672	673	674	675	676	677	678	679	680	681	682	683	684	685	686	687	688	689	690	691	692	693	694	695	696	697	698	699	700
WK-501	laua	Cugc	Agcu	Sagu	Yuuu	Qcag	Tacu	Qcag	Tacu	Naau	Sucu	Pccu	Rcgg	Rcgg	Agca	Rcgu	Sagu	Vgua	Agcu	Sagu	Qcaa	Succ	lauc	lauu	Agcc	Yuac	Tacu	Maug	Suca	Lcuu	Gggu
WK-012	laua	Cugc	Agcu	Sagu	Yuuu	Qcag	Tacu	Qcag	Tacu	Naau	Sucu	Pccu	Rcgg	Rcgg	Agca	Rcgu	Sagu	Vgua	Agcu	Sagu	Qcaa	Succ	lauc	lauu	Agcc	Yuac	Tacu	Maug	Suca	Lcuu	Gggu
WK-521	laua	Cugc	Agcu	Sagu	Yuuu	Qcag	Tacu	Qcag	Tacu	Naau	Sucu	Pccu	Rcgg	Rcgg	Agca	Rcgu	Sagu	Vgua	Agcu	Sagu	Qcaa	Succ	lauc	lauu	Agcc	Yuac	Tacu	Maug	Suca	Lcuu	Gggu
WA1-A12	laua	Cugc	Agcu	Sagu	Yuuu	Qcag	Tacu	Qcag	Tacu	Naau	Sucu	Pccu	Rcgg	Rcgg	Agca	Rcgu	Sagu	Vgua	Agcu	Sagu	Qcaa	Succ	lauc	lauu	Agcc	Yuac	Tacu	Maug	Suca	Lcuu	Gggu
WA1-F6	laua	Cugc	Agcu	Sagu	Yuuu	Qcag	Tacu	Qcag	Tacu	Naau	Sucu	Pccu	Rcgg	Rcgg	Agca	Rcgu	Sagu	Vgua	Agcu	Sagu	Qcaa	Succ	lauc	lauu	Agcc	Yuac	Tacu	Maug	Suca	Lcuu	Gggu
HU-1	laua	Cugc	Agcu	Sagu	Yuuu	Qcag	Tacu	Qcag	Tacu	Naau	Sucu	Pccu	Rcgg	Rcgg	Agca	Rcgu	Sagu	Vgua	Agcu	Sagu	Qcaa	Succ	lauc	lauu	Agcc	Yuac	Tacu	Maug	Suca	Lcuu	Gggu
RaTG13	laua	Cugc	Agcc	Sagu	Yuuu	Qcag	Tacu	Qcaa	Tacu	Naau	Suca	---	---	---	---	Rcgu	Sagu	Vgug	Agcc	Sagu	Qcaa	Suct	lauu	lauu	Agcc	Yuac	Tacu	Maug	Suca	Lcuu	Gggu

7
8
9

Fig 2. Codon and aa unified view of the codon alignments: (A) the codon alignment of HIV/SIV envelop protein gene; HV1J3-HV1MA: HIV strains, SIVCZ: chimpanzee SIV; (B) the codon alignment of the spike protein gene of coronaviruses. WK-501, WK-012, WK-521, WA1-A12, WA1-F6, HU-1: COVID-19 isolates; RaTG13: a bat coronavirus (MN996532.1) which is identified as the most recent common ancestor of COVID-19 and SARS COVs. Uppercase: amino acids; lowercase: nucleotides.

14
15