# Genome charaterization based on the Spike-614 and NS8-84 loci of SARS-CoV-2 reveals two major onsets of the COVID-19 pandemic

Xiaowen Hu[1,2#], Yaojia Mu[1#], Ruru Deng[1#], Guohui Yi[3], Lei Yao[*4], Jiaming Zhang[1*]

[1]Key Laboratory of Microbiology of Hainan Province, Institute of Tropical Bioscience and Biotechnology, Chinese Academy of Tropical Agricultural Sciences, Haikou, Hainan, 57110, China

[2]Zhanjiang Experimental Station, Chinese Academy of Tropical Agricultural Sciences, Zhanjiang, 524013, China

[3]Public Research Laboratory, Hainan Medical University, Haikou571199, China

[4]Experiment Medicine Center, The Affiliated Hospital of Southwest Medical University, 646000, Luzhou, Sichuan, China

[#]These authors contributed equally

[*] Correspondence should be addressed to J.Z. (zhangjiaming@itbb.org.cn) and L.Y. (yaolei2009@gmail.com)

Xiaowen Hu: yuyin110110110@163.com

Yaojia Mu: muyaojia@foxmail.com

Ruru Deng: 601859770@qq.com

Guohui Yi: guohuiyi6@hainmc.edu.cn

Lei Yao: yaolei2009@gmail.com

Jiaming Zhang: zhangjiaming@itbb.org.cn

# Abstract

The global COVID-19 pandemic has lasted for three years since its outbreak, however its origin is still unknown. Here, we analyzed the genotypes of 3.14 million SARS-CoV-2 genomes based on the amino acid 614 of the Spke (S) and the amino acid 48 of NS8 (nonstructural protein 8), and identified 16 linkage haplotypes. The GL haplotype (S_614G and NS8_48L) was the major haplotype driving the global pandemic and accounted for 99.2% of the sequenced genomes, while the DL haplotype (S_614D and NS8_48L) caused the pandemic in China in the spring of 2020 and accounted for approximately 60% of the genomes in China and 0.45% of the global genomes. The GS (S_614G and NS8_48S), DS (S_614D and NS8_48S) and NS (S_614N and NS8_48S) haplotypes accounted for 0.26%, 0.06%, and 0.0067% of the genomes, respectively. The main evolutionary trajectory of SARS-CoV-2 is DS→DL→GL, whereas the other haplotypes are minor byproducts in the evolution. Surprisingly, the newest haplotype GL had the oldest time of most recent common ancestor (tMRCA), which was May 1 2019 by mean, while the oldest haplotype had the newest tMRCA with a mean of October 17, indicating that the ancestral strains that gave birth to GL had been extinct and replaced by the more adapted newcomer at the place of its origin, just like the sequential rise and fall of the delta and omicron variants.  However, they arrived and evolved into toxic strains and ignited a pandemic in China where the GL strains did not exist at the end of 2019. The GL strains had spread all over the world before they were discovered, and ignited the global pandemic, which had not been noticed until the pandemic was declared in China. However, the GL haplotype had little influence in China during the early phase of the pandemic due to its late arrival as well as the strict transmission controls in China. Therefore, we propose two major onsets of the COVID-19 pandemic, one was mainly driven by the haplotype DL in China, the other was driven by the haplotype GL globally.

**Keywords:** SARS-CoV-2, virus, molecular evolution, population genetics

# Introduction

The severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) is the causal agent of COVID-19, a disease first reported in Wuhan, China [1-3]. This virus differs from all known coronaviruses,

1    including the six coronavirus species that are known to cause human disease [4], and has been

2    classified as the seventh coronavirus that can infect humans by the International Committee on

3    Taxonomy of Viruses [5]. SARS-CoV-2 quickly spread within China. A range of public health

4    interventions was used to control the epidemic, including isolation of suspected and confirmed cases,

5    transport prohibition in and out of epidemic centers, suspension of public transport, closing of

6    schools and entertainment venues, and public gathering bans and health checks [6, 7]. These

7    measures eventually brought about an end to sustained local transmission across China in April

8    2019. However, COVID-19 expands its territories, and became a global pandemic. As of 8

9    September 2021, there have been over 266 million confirmed cases and 5.2 million deaths of

10   COVID-19, reported to WHO (https://covid19.who.int/).

11       To our knowledge, the origins of SARS-CoV-2 remain elusive. Understanding how, when, and

12   where the virus was transmitted from its natural reservoir to humans is crucial for preventing future

13   coronavirus outbreaks [8]. Progresses have been made in this aspect. A bat-origin virus RaTG13

14   that has 96% sequence identities with SARS-CoV-2 has been published [9], which suggests bats as

15   a likely natural reservoir of this virus. Pangolin [10], snakes [11], turtles [12], and/or Bovidae and

16   Cricetidae [13] have been suggested to act as potential intermediate hosts that helped the virus to

17   cross the species barriers to infect humans. According to the clinical summary of the earliest cases

18   of COVID-19 (also known as 2019-nCoV), the majority of cases were exposed to the Huanan

19   Seafood Market [14], which also had wild animals. Therefore, the market was considered to be an

20   obvious candidate location of the initial zoonotic transmission event of COVID-19 [15], and both

21   epidemiological and phylogenetic approaches suggested late 2019 the occurring time of the

22   pandemic [14, 16, 17]. However, none of the wild animals from the Huanan Seafood Market were

23   tested positive for SARS-CoV-2 [18]. Some environmental samples were positive, but their viral

24   genomes were not located at the basal position of the phylogeny [19]. Moreover, some of the early

25   cases were not epidemiologically linked to the Huanan Seafood Market [16], but linked to other

26   markets [20], and animal-to-human transmission in the market has never been confirmed and should

27   not be overemphasized [21].

28       The origin and evolution of SARS-CoV-2 has been well reviewed [22, 23]. Most tracing

29   research focused on analysis of the early genomes [23-25]. For example, Tang and colleagues

1     classified the early SARS-CoV-2 genomes into two major lineages (L and S lineages), and the two

2     lineages were well defined by just two SNPs that show complete linkage across SARS-CoV-2

3     strains[26]. The L lineage constituted 70% of the sequenced genome, while the S lineage constituted

4     approximately 30% of the sequenced genome in the 103 genomes in the early phase of epidemic

5     [26]. However, higher percentage is not sufficient to substantiate a more aggressive type [27], until

6     the emergence of the delta variant [28]. Evolutionary analyses suggested the S lineage appeared to

7     be more related to coronaviruses in animals [26]. Li and colleagues outlined the early viral spread

8     in Wuhan and its transmission patterns in China and across the rest of the world [23]. Pekar and

9     colleagues analyzed the early viral genomes (before the end of April 2020) in China, and defined a

10     period between mid-October and mid-November 2019 as the plausible interval when the first case

11     of SARS-CoV-2 emerged in Hubei province, China [24]. Ruan and colleagues discovered distinct

12     set of mutations driving the waves of replacements of strains, and the split between the Asian and

13     European lines occurred before September of 2019, suggesting twin-beginning scenario of the

14     pandemic [25]. Due to the presumed undetected transmission and insufficient genome sequences in

15     the early phase of the pandemic, the origin of the virus remains uncovered. It is true that

16     coronaviruses evolve rapidly [29], and the mutation rate of SARS-CoV-2 was predicted to be

17     approximately $8.0 \times 10^{-4}$ subs/site/year [30]. However, the inheritance rate of each nucleotide should

18     be higher than 99.9992%, taken the negative selection pressure into account [26]. Upon the

19     increased coverage of the sequenced genomes in the late phase of the pandemic, we believe that the

20     secret of the viral origin may be buried in the genomes, not necessarily in the early genomes.

21        In this study, we analyzed 3.14 million genomes obtained from GISAID database, and used

22     two gene loci to anchor the haplotypes, and identified the haplotype for the most recent common

23     ancestor (MRCA) and the proximal geographical origin of SARS-CoV-2. The evolutionary

24     trajectory of the SARS-CoV-2 haplotypes is proposed.

## 25    Materials and methods

26     SARS-CoV-2 genome sequences were obtained from the GISAID database (https://www.gisaid.org)

27     on 18 October 2021. The genomes that were incomplete <29,000 nt, had low coverage with >0.5%

28     unique amino acid mutations and more than 1% 'N's were filtered. A total of 3,140,626 genome

29     sequences of SARS-CoV-2 were remained for further analysis. All the genomes were isolated from

1    human patients. To analyze the temporal-dependent evolutionary pattern of SARS-CoV-2, the

2    genome sequences were divided into groups by months according to their sample collection date,

3    and were aligned using ViralMSA [31] and Clustal-Omega [32]. The unknown and/or low quality

4    nucleotides were removed from the alignments using personalized scripts, which is available

5    (https://github.com/XiaowenH/SeqSC.git). Nucleotide diversity (Pi) and population mutation rate

6    (Theta-W) of each subgroup was calculated using Pairwise Deletion Model with DnaSP 6 [33]. The

7    Pearson correlation coefficient [34] and significant analysis was calculated by cor.test in the R

8    package (version 4.1.1) developed by the R Development Core Team (https://www.r-project.org).

9       To genotype the genomes based on the spike gene, the spike coding sequences of the 3.14

10    million genomes were aligned to the reference genome Wuhan-Hu-1 (MN908947.3), and the

11    relevant CDS was extracted and translated into peptides by using SeqKit package [35]. The peptides

12    were aligned with the reference using Clustal-Omega [32], and the site aligned with the site 614 of

13    the reference was identified using personalized script (https://github.com/XiaowenH/SeqSC.git).

14    The NS8 gene was genotyped using the same method, except that NS8 gene was used as the

15    reference.

16       To analyze the mutation accumulation pattern of each haplotype, the genome sequences were

17    divided by month, and were aligned using ViralMSA[31] and/or Clustal-Omega [32]. Phylogenetic

18    analysis was performed using Mega 11 [36]. The evolutionary network analysis was performed by

19    using the Median-joining Networks in Potpart 1.7 [37].

20       To explore the evolutionary dynamics, genomes of each haplotype were aligned with

21    ViralMSA[31] after removal of redundants by using CD-HIT [38] with a threshhold of

22    99.98%. The outputs were transformed to Nexus format by using ALTER [39]. Bayesian

23    phylodynamics analysis was performed by using BEAST 1.10 [40] with molecular clock set to strict

24    and coalescent prior set to Bayesian Skyline [24]. The posterior distribution was summarized by

25    using TRACER 1.7 [41].

26

27    **Results**

28    **Genome and mutation accumulation patterns of SARS-CoV-2**

1   A total of 3,140,626 SARS-CoV-2 genomes designated as complete and high coverage in the

2   GISAID database (www.gisaid.org) were obtained on 2021-10-18. The average length of the

3   genomes is 29,781.5 bases with the minimum and maximum lengths of 29,000 and 31,579 bases,

4   respectively. Both the sequenced genomes and confirmed cases increased exponentially during 2020,

5   and remained at high level during 2021 (Fig. 1A). Regression analysis shows that the sequenced

6   genomes and confirmed cases are positively correlated with a linear coefficient r=0.75 (p-value =

7   5.5e-05, Fig. 1B), indicating that the sequenced genomes sufficiently represented the dynamic

8   population.

9       The nucleotide diversity (Pi) was accumulated in a temporal-dependent pattern (Fig.1C,

10  supplementary Table S1). Both Pi and Theta-W (population mutation rate) increased linearly in

11  2020 with correlation coefficients r=0.99 (p=5.349e-10, Fig. 1C) and r=0.975 (p=7.331e-08, Fig.1D)

12  for Pi and Theta-W, respectively. They then varied up and down in 2021, possibly due to the shift

13  of dominant viral strains, such as the delta and the omicron variants (submitted to PlosOne).

14  **Mutations of the amino acid 84 of NS8 in the 3.14 million genomes**

15  The mutation at site 84 of the nonstructural protein ORF8  (NS8) was one of the earliest

16  sites to be used to  classified the early SARS-CoV-2 strains, and they were classified

17  into two major lineages (designated L and S) that were well defined by two different

18  SNPs (T8782C and C28144T, the latter being Orf8-S84L) [26]. The NS8 contains 121

19  amino acids, and promotes the expression of the ER unfolded protein response factor ATF6 in

20  human cells [42]. We extracted the NS8 coding DNA sequence from all the 3.14 million SARS-

21  CoV-2 genomes, and classified the genomes according to the identities at the amino acid 84. Besides

22  the two known alleles, four additional mutations were identified, including NS8_84V, NS8_84I,

23  NS8_84F, NS8_84C (Table 1). NS8_84L was the first to have been collected. Its first collection

24  date was 2019-12-24 (EPI_ISL_402123), and a total of 21 NS8_84L samples were collected in

25  Wuhan, China in December 2019, and constituted 95.5% of the early samples (Supplementary Table

26  S3). This genotype constituted 99.66% of the 3.13 million genomes that were successfully classified

27  (Table 1). NS8_84S was first collected on 2019-12-30 and constituted approximately 0.34% of the

28  total genomes. The numbers of the newly identified genomes were very few (Table 1).

1    Approximately 0.4% genomes were unidentified due to low quality in the region (Table 1).

2

3

4    **Table 1.** Alleles of the amino acid 84 in the nonstructural protein 8 (NS8).

| Allele | Sequence (75-93) | Genomes | Percentage | First collect | Collection place |
|---|---|---|---|---|---|
| NS8_TG13 | DIGNYTVSC**S**PFTINCQEP | RaTG13 ref. | | 2013-07-24 | Yunan,China |
| NS8_wiv04 | DIGNYTVSC**L**PFTINCQEP | GISAID ref. | | 2019-12-30 | Wuhan,China |
| NS8_84L | DIGNYTVSC**L**PFTINCQEP | 3,115,900 | 99.6564* | 2019-12-24 | Wuhan,China |
| NS8_84S | DIGNYTVSC**S**PFTINCQEP | 10,732 | 0.3432 | 2019-12-30 | Wuhan,China |
| NS8_84V | DIGNYTVSC**V**PFTINCQEP | 3 | 0.0001 | 2020-09-15 | South Korea |
| NS8_84I | DIGNYTVSC**I**PFTINCQEP | 4 | 0.0001 | 2021-05-05 | USA |
| NS8_84F | DIGNYTVSC**F**PFTINCQEP | 2 | 0.00006 | 2020-02-02 | Sichuan,China |
| NS8_84C | DIGNYTVSC**C**PFTSNCQEP | 1 | 0.00003 | 2020-03-12 | USA |
| unknown | | 13,985 | | | |
| Total | | 3,140,626 | | | |

5    *Note: Percentage in the genomes that have been successfully classified.

6

7    **Genome and mutation accumulation and phylogenetic analysis of NS8 gene**

8    NS8_84L was dominant over NS8_84S all the time (Supplementary Table S3). The genome number

9    of NS8_84L increased exponentially from December 2019 to March 2021, and kept at high numbers

10    (Fig. 2A), while NS8_84S had a peak in March 2020, and remained at low numbers the rest of time

11    (Fig. 2B). Even in its peak month, its genome number was only 12% of NS8_84L (supplementary

12    Table S3). The genome number of NS8_84L is positively correlated with the number of the

13    confirmed cases with a coefficient r=0.88 (p=8.914e-05, Fig. 2C), while the genome number of

14    NS8_84S is not (r=-0.32, p=0.29, Fig.2D). Therefore, NS8_84L was the main genotype driving the

15    growth of confirmed cases worldwide.

16        The nucleotide diversity of both genotypes was accumulated linearly in 2020. Regression

17    analysis revealed first occurrence dates of approximately 14 October 2019 (95% confidence interval:

18    18 September to 9 November) for the genotype NS8_84L (Fig. 2E), and 28 September 2019 (95%

7

1    confidence interval: 6 August to 21 November) for  genotype and NS8_84S (Fig.2F), respectively.

2    Therefore, NS8_84S occurred earlier than NS8_84L. Phylogenetic analysis using both peptide and

3    RNA sequences of the NS8 showed that NS8_84S was closest to the root (Fig. 2G, H), and the

4    phylogenies by using different methods all suggested an earlier occurrence of NS8_84S than

5    NS8_84L consistent (Supplementary Figure S2 and S3), which was supported by that the bat SARS-

6    CoV related viruses all have NS8_84S [26].

7        Taken together, NS8_84S is an earlier genotype than NS8_84L, and the latter may be

8    originated from mutation of NS8_84S

9

10   **Mutations of the amino acid 614 of the spike protein in the 3.14 million genomes**

11   Mutation of D614G in Spike protein was the major driving force in the COVID-19 pandemic, the

12   virulent strains alpha, delta, and omicron all carried this mutation. The Spike (S) is a surface

13   projection glycoprotein. Mutations in this protein have previously been associated with altered

14   pathogenesis and virulence in other coronaviruses [43]. We extracted the S gene sequences from

15   the 3.14 genomes, and classified the genomes based on mutations at the amino acid 614. Seven

16   alleles were identified, including S_614D, S_614G, S_614N, S_614V, S_614A, S_614V, S_614C

17   (Table 2).

18       The S_614D was the earliest to have been collected, and was first collected in Wuhan, China

19   on 2019-12-24 (EPI_ISL_402123, Table 2). In outside China, it was first collected in Thailand on

20   2020-01-08, Nepal on 2020-01-13, USA on 2020-01-19, and France on 2020-01-23. This genotype

21   constituted the second largest number of the genomes with a percentage of 0.7%.

22       The S_614G accounted for 99.69% of the genomes (Table 2), and was first collected on 2020-

23   01-01 in Argentina (EPI_ISL_4405694), followed by nine samples collected on 2020-01-03 in USA

24   (EPI_ISL_3537067,        EPI_ISL_3537066,        EPI_ISL_3537065,        EPI_ISL_3537064,

25   EPI_ISL_3537063,        EPI_ISL_3537062,        EPI_ISL_3537061,        EPI_ISL_3537060,

26   EPI_ISL_3537059). It was first collected in Australia (EPI_ISL_3568416) on 2020-01-08, in Japan

27   (EPI_ISL_2671842) on 2020-01-09, in India on 2020-01-12, and in Africa (EPI_ISL_2716636) on

28   2020-01-14. This variant was collected in Zhejiang (EPI_ISL_422425) and Sichuan

29   (EPI_ISL_451345) provinces of China on 2020-01-24.

1    The S_614N had the third largest number with 200 genomes. It was first collected in England

2    on 2020-03-24. The S_614S was characterized by 61 genomes, while the other genotypes had very

3    few numbers (Table 2).

4    **Table 2.** Alleles of the amino acid 614 in the spike protein (spike_614) of SARS-CoV-2.

| Allele | Sequence (605-623) | Genomes | Percentage | First collect | Collection place |
|--------|--------------------|---------|------------|---------------|------------------|
| S_TG13 | SNQVAVLYQ**D**VNCTEVPVA | reference | | 2013-07-24 | Yunnan, China |
| S_wiv04 | SNQVAVLYQ**D**VNCTEVPVA | reference | | 2019-12-30 | Wuhan, China |
| S_614G | SNQVAVLYQ**G**VNCTEVPVA | 2,846,092 | 99.293 | 2020-01-01 | France, Argentina |
| S_614D | SNQVAVLYQ**D**VNCTEVPVA | 19,992 | 0.697 | 2019-12-24 | Wuhan, China |
| S_614N | SNQVAVLYQ**N**VNCTEVPVA | 200 | 0.00698 | 2020-03-24 | England |
| S_614S | SNQVAVLYQ**S**VNCTEVPVA | 61 | 0.00213 | 2020-05-21 | USA |
| S_614A | SNQVAVLYQ**A**VNCTEVPVA | 11 | 0.00038 | 2020-07-13 | South Korea |
| S_614V | SNQVAVLYQ**V**VNCTEVPVA | 2 | 0.00007 | 2020-07-10 | USA |
| S_614C | SNQVAVLYQ**C**VNCTEVPVA | 1 | 0.00003 | 2021-07-30 | Cambodia |
| unknown | | 274,267 | | | |
| Total | | 3,140,626 | | | |

5

6

7    **Mutation accumulation pattern and phylogenetic analysis of Spike variants**

8    Similar to NS8 gene, S_614G and S_614D constituted 99.99% of the sequenced genomes

9    (Supplementary Table S4). At the same time, S_614G was dominant over S_614D all the time. The

10   genome number of S_614G increased exponentially from December 2019 to March 2021, and kept

11   at high numbers (Fig.3A), while S_614D reached a peak in March 2020, and remained at low

12   numbers for the rest of time (Fig. 3B，supplementary Table S4). Regression analysis shows that

13   the genome number of S_614G is positively correlated with the number of confirmed cases globally

14   with a coefficient r= 0.69 (p= 0.0031, Fig.3C), while the genome number of S_614D had no

15   correlation (r= -0.39, p=0.14, Fig.3D). Therefore, S_614G was the main haplotype driving the

16   growth of the confirmed cases worldwide.

17   The nucleotide diversity of S_614G genomes was accumulated linearly in 2020 with

1   $R^2$=0.9733 (Fig. 3E), while the nucleotide diversity of S_614D increased linearly from December

2   2019 to August 2020, but decreased sharply in October 2020 (Fig.3F), which was probably resulted

3   by the strict lockdown in China, since this variant mainly present in China (see below).

4       Phylogenetic analysis using genome sequences of S_614 alleles indicated that S_614N and

5   S_614S were the closest relatives to the root in the ML tree (Fig. 3G), however, consistent results

6   were not obtained by using the NJ and ME methods (supplementary Figure S4), and/or by the

7   phylogenenies of peptide sequences (Fig. 3H, supplementary Figure S5). Since S_614S is used in

8   the bat spikes, S_614S should be the ancestral allele, and S_614N belongs to an allele that evolved

9   early in evolutionary history.

10

11  **Genome classification based on the two loci of NS8 and S**

12  The six NS8_84 alleles and the seven spike_614 alleles should theoretically form 42 linkage types.

13  However, only 16 haplotypes were identified in the 3.14 million genomes (Table 3). These

14  haplotypes are termed by two capital letters, the first letter represents the amino acid 614 of spike,

15  while the second letter represents the amino acid 84 of NS8. For example, S_614G+NS8_84L is

16  termed GL, S_614G+NS8_84S is termed GS. Four main haplotypes GL, GS, DL, and DS accounted

17  for 99.99% of the total classified genomes. The other haplotypes included GV, GI, GF, DF, DC,

18  NL, NS, SL, AL, AS, VL, CL (Table 3).

19      GL haplotype was dominant over the other main haplotypes with 2.8 million sequenced

20  genomes (99.24% of all genomes). It was dominant all the time since March 2020 with hundreds of

21  thousands of genomes (Figure 4A). GL was first collected in Argentina (EPI_ISL_4405694) on

22  2020-01-01 (submission dates are 2021-09-22), and was collected in USA (EPI_ISL_3537059,

23  EPI_ISL_3537060, EPI_ISL_3537061, EPI_ISL_3537062) on 2020-01-03. This haplotype was

24  first collected in Sichuan (EPI_ISL_451345) and Zhejiang (EPI_ISL_422425) China on 2020-01-

25  24.

26      GS haplotype was first collected in Wuhan, China (EPI_ISL_412982) on 2020-02-07, which

27  was the only sample collected globally in February 2020. It did not show up in China again, due to

28  its rareness. Genome analysis showed that this genotype came from a recombination between GL

29  and DS haplotype, since the S variant came from a variation of C28144T, which is closely linked

10

1  with a variation of C8782T [26], however, the variation C28144T in the GS haplotype is not linked

2  with C8782T. The GS haplotype showed up again in Spain on 2020-03-03, and a total of 19 GS

3  genomes were identified in 2020-03 globally, including 10 in Spain, seven in USA, one in Scotland,

4  and one in Belgium. Its number increased to a few hundreds each month in 2021 (Fig. 4B). All these

5  strains did not have the linked variation C8782T, therefore, they may have all come from

6  recombination between GL and DS haplotypes, or reverse mutations.

7  The DS and DL genomes reached more than a thousand only in March and April 2020 (Fig.

8  4C), and kept at very low numbers after May 2020, possibly due to the strict lockdown policy in

9  mainland China, since DL and DS contributed 92.6% of the genomes in mainland China before June

10  2020, but contributed very few proportions outside China (Fig. 4D).

11  The genome numbers and the first occurrence of other genotypes are provided in Table 3.

12  Regression analysis showed that GL was the major haplotype driving the increase of confirmed

13  cases worldwide with a coefficient r=0.77, p= 2.571e-05 (Fig. 4E), while GS was weak (r=0.53, p=

14  0.01152, Fig. 4F). The genome numbers of DL and DS were negatively correlated with the

15  confirmed cases, however, not significant at 5% significance level (r=-0.4, p>0.05, Fig. 4G, H).

16  **Table 3.** Genome numbers of Spike_614 and NS8_84 mutation linkage types and their first collection date.

| | NS8_84L | NS8_84S | NS8_84V | NS8_84I | NS8_84F | NS8_84C | total |
|---|---|---|---|---|---|---|---|
| S_614G | 2,831,567 | 1,622 | 3 | 4 | 1 | 0 | 2,833,197 |
| | 2020-01-01 | 2020-02-07 | 2020-09-15 | 2021-05-05 | 2021-05-06 | | |
| | Argentina | Wuhan,China | South Korea | USA | Sweden | | |
| S_614D | 12,865 | 7,016 | 0 | 0 | 1 | 1 | 19,883 |
| | 2019-12-24 | 2019-12-30 | | | 2020-02-02 | 2020-03-12 | |
| | Wuhan,China | Wuhan,China | | | Sichuan,China | USA | |
| S_614N | 9 | 191 | 0 | 0 | 0 | 0 | 200 |
| | 2020-03-24 | 2020-06-10 | | | | | |
| | England | Mali | | | | | |
| S_614S | 61 | 0 | 0 | 0 | 0 | 0 | 61 |
| | 2020-05-21 | | | | | | |
| | USA | | | | | | |
| S_614A | 4 | 7 | 0 | 0 | 0 | 0 | 11 |
| | 2020-12-12 | 2020-07-13 | | | | | |
| | USA | South Korea | | | | | |
| S_614V | 2 | 0 | 0 | 0 | 0 | 0 | 2 |
| | 2020-07-10 | | | | | | |
| | USA | | | | | | |
| S_614C | 1 | 0 | 0 | 0 | 0 | 0 | 1 |
| | 2021-07-30 | | | | | | |
| | Cambodia | | | | | | |
| total | 2,844,509 | 8,836 | 3 | 4 | 2 | 1 | 2,853,355 |

17

18  **Diversity dynamics of the main haplotypes**

19  Mutation accumulation curve showed that the nucleotide diversity of haplotypes GL, DL, and DS

20  all had a period of linear increase. The difference was that GL had a high initial Pi in January 2020

21  followed by a sharp decrease in February 2020 before the linear increase (Figure 4I). The

1    nonsynonymous and synonymous mutation ratios (Ka/Ks) are between 0.65 to 0.29 (Supplementary

2    Fig.S6), indicating the presence of a strong purify selection, and many strains may have failed in

3    further transmission till an adaptive strain emerged, which was also pointed out by Pekar and

4    colleagues [24].

5    In contrast to GL haplotype, DL and DS haplotypes went straight to the linear period without

6    initial high diversity, with correlation coefficients 0.995 (p= 3.723e-06) and 0.961 (p= 1.461e-5) for

7    DL (Fig. 4J) and DS (Fig. 4K), respectively. These results suggested that DL and DS may have gone

8    through the selection period when SARS-CoV-2 was discovered. However, the linear accumulation

9    periods for DL and DS haplotypes were relatively shorter compared to that of GL, possibly due to

10    the strong transmission control measures carried out in China, since these two haplotypes were the

11    majority haplotypes in China (Fig. 4D).

12    The mutation accumulation curve of GS did not show any linear period, but its diversity was

13    mostly high (Fig. 4L), which was even the highest among the main haplotypes most of the months

14    in 2020 (Supplementary Figure S7). These results suggested that GS may have evolved many times

15    by recombination or reverse mutation, and had never obtained strong transmission ability in the

16    human population for some reason, which was also supported by the lowest number of the

17    sequenced genomes in the four main haplotypes (Fig. 4B, supplementary Table S5).

18

19    **Phylogenetic and phylodynamics analysis of the haplotypes**

20    Among the three main haplotypes (GL, DL, and DS), DS is the closest haplotype to bat genomes,

21    since most bat genomes have Spike-614D and NS8-84S [25]. The evolutionary trajectory of SARS-

22    CoV-2 should be DS → DL → GL. The rest haplotypes are minor haplotypes resulting from

23    recombination, reverse mutation, and/or mutation in the adjacent sites. Phylogenetic analysis

24    revealed that the NS haplotype together with AS and DS are the closest to the root by using all the

25    ML, MP, NJ, and ME methods (Fig. 5, supplementary Figure S8), followed DL and GL. Therefore,

26    the NS, AS, and DC may have arrived from mutations of DS. The GS haplotype are distributed

27    widely in the DS, DL, and GL haplotypes, indicating their multiple origins (recombination between

28    GL and DS, or mutation from either DS or GL), which is in agreement with previous observation

29    of its always high Pi (Fig. 4L). NL and DF haplotypes are located in the clade of DL, indicating its

1  origin from DL, whereas GF, GV, SL, VL, and AL were derived from GL.

2  Bayesian phylodynamics analysis [40] was performed to explored the evolutionary dynamics

3  of the major haplotypes by using a Bayesian Skyline approach. Surprisingly, the newest haplotype

4  GL in the main haplotypes was inferred to have the oldest tMRCA with a mean of May 1 (95%

5  HPD interval: Feb 8 to Aug 4 2019), whereas the oldest haplotype DS had the newest tMRCA with

6  a mean of October 17 2019 (95% HPD: September 11 to November 23 2019), while the tMRCA of

7  DL was inferred to locate in September with a mean of September 20 (95% HPD: August 11 to

8  October 24 2019). The tMRCA of the three haplotypes had a mean of February 18 (95% HPD: Nov

9  5 2018 to Jun 4 2019). These results indicate that the common ancestor of the three haplotypes may

10  have evolved before February 2019, and the three haplotypes co-existed before April 2019. The

11  ancestor strains of DS and DL that gave birth to GL have been extinct upon the outbreak of COVID-

12  19, therefore, much recent tMRCAs were estimated by using the extant genome sequences.

13  **Proposed evolutionary trajectory of SARS-CoV-2**

14  Based on the results of mutation accumulation curves, phylogenetic and phylodynamics analysis,

15  and combined with codon usage of the mutants by single nucleotide mutation principle, we propose

16  a possible evolutionary trajectory of SARS-CoV-2 haplotypes as shown in Figure 7.

17  The DS haplotypes was the ancestral haplotype of SARS-CoV-2, and may have occurred in

18  January 2019, and evolved to DL in February 2019, and further evolved to GL in April 2019. The

19  ancestral strains that gave birth to GL went distinct and replaced by the more adapted newcomer at

20  the place of its origin. However, they arrived and evolved into toxic strains in China where the GL

21  strains did not exist, and ignited the COVID-19 pandemic in China at the end of 2019. The GL

22  strains had spread all over the world before the pandemic began in China, and ignited the global

23  pandemic, which had not been noticed until it was reported in China.

24

25  # Discussion

26  Tracing the origin of SARS-CoV-2 is critical for preventing future spillover of the virus [44]. This

27  is a routine process in infectious disease prevention and control. Despite of tremendous efforts, the

28  origin of SARS-CoV-2 remains unclear. Bats have been recognized as the natural reservoirs of a

29  large variety of viruses [45]. SARS-CoV-1 that caused a pandemic in China in 2003 and Middle

1  East Respiratory Syndrome Coronavirus (MERS-CoV) are both suggested to be originated from

2  bats [46-48]. A bat-origin coronavirus RaTG13 has the most similar genome compared to SARS-

3  CoV-2 with 96.2% identities on whole genome level [49]. Several other bat-origin coronavirus with

4  highly similar genomic sequences compared to SARS-CoV-2 have also been found in different

5  countries [50-52]. Although the receptor binding domains (RBD) of the spike proteins of RaTG13

6  was only 89.2% compared to SARS-CoV-2, it could bind to human ACE2 (hACE2), and RaTG13

7  pseudovirus could transduce cells expressing hACE2 with low efficiency [53]. Conversely, the

8  SARS-CoV-2 spike protein RBD could bind to bACE2 from *Rhinolophus macrotis* with

9  substantially lower affinity compared with that to hACE2, and its infectivity to host cells expressing

10  bACE2 was confirmed with pseudotyped SARS-CoV-2 virus and SARS-CoV-2 wild virus [54].

11  Based on these facts, SARS-CoV-2 could have more likely originated from bats. However,

12  intermediate hosts are needed for bat-origin coronaviruses to acquire sufficient mutations so as

13  to infect humans [55]. Two SARS-CoV-2-related coronaviruses with RBD highly similar to that of

14  SARS-CoV-2 were found in Malayan pangolins (*Manis javanica*)[10, 56]. However, their overall

15  genomic similarity compared to SARS-CoV-2 was both low (<93%), which suggested that

16  pangolins were unlikely to be the intermediate host for SARS-CoV-2. So far, the direct evolutionary

17  progenitor of SARS-CoV-2 remains unclear, whether bats are the original reservoir hosts, how,

18  when and where SARS-CoV-2 was transmitted from animals to humans remain mysteries.

19  In this study, we used two gene loci to classify the haplotypes of 3.14 million SARS-CoV-2

20  genomes, and identified seven S_614 alleles and six NS8_48 alleles, and 16 linkage types.

21  Phylogenetic and phylodynamics, mutation accumulation curve and codon usage analysis proposed

22  the evolutionary trajectory of the 16 haplotypes (Fig. 7). The DS haplotype was proposed to be the

23  oldest haplotype and probably represents the haplotype of the most recent common ancestor (MRCA)

24  of all SARS-CoV-2 strains. However, the ancestral strain of the DS haplotype was not found due to

25  its low adaptability compared to GL haplotype. However, its possible occurring time was estimated

26  by using the currently available genomes of the three main haplotypes, and the mean tMRCA was

27  estimated to be February 18 2019. Although the estimated tMRCAs of SARS-CoV-2 found in

28  humans means the time to the most recent common ancestor of the viral variants, and should not be

29  interpreted as the timing of the viral jump from animal hosts to humans, the tMRCA can be regarded

as a most recent time, and the actual occurring time may be much earlier. Imagine that an ancestral SARS-CoV-2 invaded human populations, e.g. 20 years ago. Since then, the viral population may have undergone a series of strain replacements due to genetic drift and selective sweep, just like the sequential rise and fall of the delta and omicron variants. The continual losses and gains of genetic diversity may result in the re-building of the extant diversity, which may result in a very recent tMRCA. Anyway, the actual time of emergence should be much earlier than the estimated tMRCA based on the extant genome diversity.

The evolutionary trajectory we proposed suggested two major onsets of the COVID-19 pandemic (Fig. 7), one was in China mainly driven by the haplotype DL, the other was driven by the haplotype GL outside China. Since GL had already spread all over the world before the pandemic was delcared, its place of origin is not known, however, this haplotype was suggested to have emerged in Europe in a recent report by Wu and colleagues [25]. The GL haplotype had little influence in China during the early phase of the pandemic due to its late arrival and the strict transmission controls in China.

# Acknowledgements

**Author contributions:** Jiaming Zhang and Lei Yao conceived and designed the experiments. Yaojia Mu, Xiaowen Hu, Ruru Deng, Xuepiao Sun, Guohui Yi, and Lei Yao performed the analysis. Xiaowen Hu wrote the personalized scripts. Jiaming Zhang wrote the draft. All authors revised and approved the manuscript.

**Conflicts of interest statement:** The authors declare no conflict of interest.

# Supplementary data

1    **Figure S1.** Temporal curve of Theta-W (population mutation rate) of the SARS-CoV-2 globally.

2    **Figure S2.** Molecular phylogenetic trees inferred by using peptide sequences of NS8; A, NJ tree,

3    B, ME tree, C, MP tree, D, ML tree. All ambiguous positions were removed for each sequence

4    pair in NJ and ME tree, or positions with less than 95% site coverage were eliminated in MP and

5    ML tree. The percentage in 1000 replicates of trees in which the associated taxa clustered together

6    is shown next to the branches. There were a total of 121 positions (NJ and ME) or 118 positions

7    (MP and ML) in the final dataset. Evolutionary analyses were conducted in MEGA11[36].

8    **Figure S3.** Molecular phylogenetic trees inferred by using RNA sequences of NS8; A, NJ tree, B,

9    ME tree, C, MP tree, D, ML tree. Codon positions included were 1st+2nd+3rd. All ambiguous

10   positions were removed for each sequence pair in NJ and ME tree, or positions with less than 95%

11   site coverage were eliminated in MP and ML tree. The percentage in 1000 replicates of trees in

12   which the associated taxa clustered together is shown next to the branches. There were a total of

13   366 positions (NJ and ME) or 354 positions (MP and ML) in the final dataset.

14   **Figure S4.** Figure S4. Molecular phylogenetic trees inferred by using RNA sequences of Spike;

15   A, NJ tree, B, ME tree, C, MP tree, D, ML tree. Codon positions included were 1st+2nd+3rd. All

16   ambiguous positions were removed for each sequence pair in NJ and ME tree, or positions with

17   less than 95% site coverage were eliminated in MP and ML tree. The percentage in 1000

18   replicates of trees in which the associated taxa clustered together is shown next to the branches.

19   **Figure S5.** Molecular phylogenetic trees inferred by using peptide sequences of Spike; A, NJ

20   tree, B, ME tree, C, MP tree, D, ML tree. All ambiguous positions were removed for each

21   sequence pair in NJ and ME tree, or positions with less than 95% site coverage were eliminated in

22   MP and ML tree. The percentage in 1000 replicates of trees in which the associated taxa clustered

23   together is shown next to the branches. There were a total of 121 positions (NJ and ME) or 118

24   positions (MP and ML) in the final dataset.

25   **Figure S6.** Accumulation curves of nonsynonymous (Ka) and synonymous (Ks) mutations. A, Ka

26   and Ks mutation rates; B, Ka/Ks ratios.

27   **Figure S7.** Temporal-dependent diversity (Pi) of four main genotypes.

28   **Figure S8.** Molecular phylogenetic trees inferred by using the genome sequence; A, ME tree, B,

29   MP tree. All ambiguous positions were removed for each sequence pair in ME tree, or positions

1   with less than 95% site coverage were eliminated in MP tree. The percentage in 1000 replicates of

2   trees in which the associated taxa clustered together is shown next to the branches.

3   **Table S1.** Temporal dependent nucleotide diversity (Pi) of genomes by collection date.

4   **Table S2.** Uneven distribution of diversity in the SARS-CoV-2 genome.

5   **Table S3.** Genome mutation accumulation of main NS8_84 genotypes.

6   **Table S4.** Genome and mutation accumulation of main spike genotypes to the GISAID database.

7   **Table S5.** Genome and mutation accumulation of main spike genotypes.

8   **Table S6.** Distribution of NS genotypes worldwide

9

# References

11  1.  Zhu N, Zhang D, Wang W, Li X, Yang B, Song J, et al. A novel coronavirus from patients with
12      pneumonia in China, 2019. N Engl J Med. 2020;382(8):727-33. doi: 10.1056/NEJMoa2001017.
13      PubMed PMID: 31978945.

14  2.  Working Group of Novel Coronavirus PUMCH. Diagnosis and clinical management of 2019 novel
15      coronavirus infection: an operational recommendation of Peking Union Medical College Hospital
16      (V2.0). Zhonghua Nei Ke Za Zhi. 2020;59(3):186-8. doi: 10.3760/cma.j.issn.0578-1426.2020.03.003.
17      PubMed PMID: 32023681.

18  3.  Wu F, Zhao S, Yu B, Chen YM, Wang W, Song ZG, et al. A new coronavirus associated with human
19      respiratory disease in China. Nature. 2020;579(7798):265-9. doi: 10.1038/s41586-020-2008-3.
20      PubMed PMID: 32015508.

21  4.  Su S, Wong G, Shi W, Liu J, Lai ACK, Zhou J, et al. Epidemiology, genetic recombination, and
22      pathogenesis of coronaviruses. Trends Microbiol. 2016;24(6):490-502.

23  5.  Coronaviridae Study Group of the International Committee on Taxonomy of V. The species Severe
24      acute respiratory syndrome-related coronavirus: classifying 2019-nCoV and naming it SARS-CoV-
25      2. Nat Microbiol. 2020;5(4):536-44. doi: 10.1038/s41564-020-0695-z. PubMed PMID: 32123347.

26  6.  Tian H, Liu Y, Li Y, Wu C-H, Chen B, Kraemer MUG, et al. An investigation of transmission control
27      measures during the first 50 days of the COVID-19 epidemic in China. Science.
28      2020;368(6491):638-42. doi: DOI: 10.1126/science.abb6105.

29  7.  Chen S, Yang J, Yang W, Wang C, Barnighausen T. COVID-19 control in China during mass
30      population movements at New Year. Lancet. 2020;395(10226):764-6. Epub 2020/02/28. doi:
31      10.1016/S0140-6736(20)30421-9. PubMed PMID: 32105609; PubMed Central PMCID:
32      PMCPMC7159085.

33  8.  Wang Q, Chen H, Shi Y, Hughes AC, Liu WJ, Jiang J, et al. Tracing the origins of SARS-CoV-2: lessons
34      learned from the past. Cell Res. 2021;31(11):1139-41. Epub 2021/10/01. doi: 10.1038/s41422-021-

1          00575-w. PubMed PMID: 34588626; PubMed Central PMCID: PMCPMC8480455.

2    9.   Zhou P, Yang X-L, Wang X-G, Hu B, Zhang L, Zhang W, et al. A pneumonia outbreak associated with
3          a new coronavirus of probable bat origin. Nature. 2020;579:270-3.

4   10.  Zhang T, Wu Q, Zhang Z. Probable pangolin origin of SARS-CoV-2 associated with the COVID-19
5          outbreak. Curr Biol. 2020. doi: 10.1016/j.cub.2020.03.022. PubMed PMID: 32197085.

6   11.  Ji W, Wang W, Zhao X, Zai J, Li X. Cross-species transmission of the newly identified coronavirus
7          2019-nCoV. J Med Virol. 2020;92(4):433-40. doi: 10.1002/jmv.25682. PubMed PMID: 31967321;
8          PubMed Central PMCID: PMCPMC7138088.

9   12.  Liu Z, Xiao X, Wei X, Li J, Yang J, Tan H, et al. Composition and divergence of coronavirus spike
10         proteins and host ACE2 receptors predict potential intermediate hosts of SARS-CoV-2. J Med Virol.
11         2020; 92(6):595-601. doi: 10.1002/jmv.25726. PubMed PMID: 32100877.

12   13.  Luan J, Jin X, Lu Y, Zhang L. SARS-CoV-2 spike protein favors ACE2 from Bovidae and Cricetidae. J
13         Med Virol. 2020;92(9):1649-56. doi: 10.1002/jmv.25817. PubMed PMID: 32239522; PubMed
14         Central PMCID: PMCPMC7228376.

15   14.  Huang C, Wang Y, Li X, Ren L, Zhao J, Hu Y, et al. Clinical features of patients infected with 2019
16         novel coronavirus in Wuhan, China. Lancet. 2020;395(10223):497-506. doi: 10.1016/S0140-
17         6736(20)30183-5. PubMed PMID: 31986264.

18   15.  Lu R, Zhao X, Li J, Niu P, Yang B, Wu H, et al. Genomic characterisation and epidemiology of 2019
19         novel coronavirus: implications for virus origins and receptor binding. Lancet.
20         2020;395(10224):565-74. doi: 10.1016/S0140-6736(20)30251-8. PubMed PMID: 32007145.

21   16.  Li Q, Guan X, Wu P, Wang X, Zhou L, Tong Y, et al. Early transmission dynamics in Wuhan, China,
22         of novel coronavirus-infected pneumonia. N Engl J Med. 2020;382(13):1199-207. Epub
23         2020/01/30. doi: 10.1056/NEJMoa2001316. PubMed PMID: 31995857; PubMed Central PMCID:
24         PMCPMC7121484.

25   17.  Zhang X, Tan Y, Ling Y, Lu G, Liu F, Yi Z, et al. Viral and host factors related to the clinical outcome
26         of COVID-19. Nature. 2020;583(7816):437-40. Epub 2020/05/21. doi: 10.1038/s41586-020-2355-
27         0. PubMed PMID: 32434211.

28   18.  Wang H, Zhao W. WHO-Convened Global Study of Origins of SARS-CoV-2: China Part (Text Extract).
29         Infectious Diseases & Immunity. 2021;1(3):125-32. doi: doi: 10.1097/ID9.0000000000000017.

30   19.  Hill V, Rambaut A. Phylodynamic analysis of SARS-CoV-2 | Update 2020-03-06. 2020.

31   20.  Holmes EC, Goldstein SA, Rasmussen AL, Robertson DL, Crits-Christoph A, Wertheim JO, et al. The
32         origins of SARS-CoV-2: A critical review. Cell. 2021;184(19):4848-56. Epub 2021/09/05. doi:
33         10.1016/j.cell.2021.08.017. PubMed PMID: 34480864; PubMed Central PMCID: PMCPMC8373617.

34   21.  Nishiura H, Linton NM, Akhmetzhanov AR. Initial cluster of novel coronavirus (2019-nCoV)
35         infections in Wuhan, China is consistent with substantial human-to-human transmission. J Clin
36         Med. 2020;9(2):488;doi:10.3390/jcm9020488. doi: 10.3390/jcm9020488. PubMed PMID:
37         32054045.

38   22.  Cui J, Li F, Shi ZL. Origin and evolution of pathogenic coronaviruses. Nat Rev Microbiol.

2019;17(3):181-92. Epub 2018/12/12. doi: 10.1038/s41579-018-0118-9. PubMed PMID: 30531947; PubMed Central PMCID: PMCPMC7097006.

23. Li J, Lai S, Gao GF, Shi W. The emergence, genomic diversity and global spread of SARS-CoV-2. Nature. 2021;600(7889):408-18. Epub 2021/12/10. doi: 10.1038/s41586-021-04188-6. PubMed PMID: 34880490.

24. Pekar J, Worobey M, Moshiri N, Scheffler K, Wertheim JO. Timing the SARS-CoV-2 index case in Hubei province. Science. 2021;372(6540):412-7. Epub 2021/03/20. doi: 10.1126/science.abf8003. PubMed PMID: 33737402; PubMed Central PMCID: PMCPMC8139421.

25. Ruan Y, Wen H, Hou M, He Z, Lu X, Xue Y, et al. The twin-beginnings of COVID-19 in Asia and Europe – One prevails quickly. National Science Review. 2021. doi: 10.1093/nsr/nwab223.

26. Tang X, Wu C, Li X, Song Y, Yao X, Wu X, et al. On the origin and continuing evolution of SARS-CoV-2. Natl Sci Rev. 2020;7(6):1012-23. Epub 2020/06/01. doi: 10.1093/nsr/nwaa036. PubMed PMID: 34676127; PubMed Central PMCID: PMCPMC7107875.

27. MacLean OA, Orton RJ, Singer JB, Robertson DL. No evidence for distinct types in the evolution of SARS-CoV-2. Virus Evolution. 2020;6(1). doi: 10.1093/ve/veaa034.

28. Davis C, Logan N, Tyson G, Orton R, Harvey WT, Perkins JS, et al. Reduced neutralisation of the Delta (B.1.617.2) SARS-CoV-2 variant of concern following vaccination. PLoS Pathog. 2021;17(12):e1010022. Epub 2021/12/03. doi: 10.1371/journal.ppat.1010022. PubMed PMID: 34855916; PubMed Central PMCID: PMCPMC8639073.

29. Hanada K, Suzuki Y, Gojobori T. A large variation in the rates of synonymous substitution for RNA viruses and its relationship to a diversity of viral infection and transmission modes. Mol Biol Evol. 2004;21(6):1074-80. doi: 10.1093/molbev/msh109. PubMed PMID: 15014142; PubMed Central PMCID: PMCPMC7107514.

30. Lai A, Bergna A, Acciarri C, Galli M, Zehender G. Early phylogenetic estimate of the effective reproduction number of SARS-CoV-2. J Med Virol. 2020:DOI:10.1002/jmv.25723. doi: 10.1002/jmv.25723. PubMed PMID: 32096566.

31. Moshiri N. ViralMSA: massively scalable reference-guided multiple sequence alignment of viral genomes. Bioinformatics (Oxford, England). 2021;37(5):714-6. Epub 2020/08/21. doi: 10.1093/bioinformatics/btaa743. PubMed PMID: 32814953.

32. Sievers F, Higgins DG. Clustal Omega for making accurate alignments of many protein sequences. Protein Sci. 2018;27(1):135-45. doi: 10.1002/pro.3290. PubMed PMID: 28884485; PubMed Central PMCID: PMCPMC5734385.

33. Rozas J, Ferrer-Mata A, Sanchez-DelBarrio JC, Guirao-Rico S, Librado P, Ramos-Onsins SE, et al. DnaSP 6: DNA sequence polymorphism analysis of large data sets. Mol Biol Evol. 2017;34(12):3299-302. doi: 10.1093/molbev/msx248. PubMed PMID: 29029172.

34. Pearson WH. Estimation of a correlation coefficient from an uncertainty measure. Psychometrika. 1966;31(3):421-33. Epub 1966/09/01. doi: 10.1007/BF02289473. PubMed PMID: 5221136.

35. Shen W, Le S, Li Y, Hu F. SeqKit: A Cross-Platform and Ultrafast Toolkit for FASTA/Q File

Manipulation. PLoS One. 2016;11(10):e0163962. Epub 2016/10/06. doi: 10.1371/journal.pone.0163962. PubMed PMID: 27706213; PubMed Central PMCID: PMCPMC5051824.

36. Tamura K, Stecher G, Kumar S. MEGA11: Molecular Evolutionary Genetics Analysis Version 11. Mol Biol Evol. 2021;38(7):3022-7. Epub 2021/04/24. doi: 10.1093/molbev/msab120. PubMed PMID: 33892491; PubMed Central PMCID: PMCPMC8233496.

37. Bandelt HJ, Forster P, Rohl A. Median-joining networks for inferring intraspecific phylogenies. Mol Biol Evol. 1999;16(1):37-48. Epub 1999/05/20. doi: 10.1093/oxfordjournals.molbev.a026036. PubMed PMID: 10331250.

38. Fu L, Niu B, Zhu Z, Wu S, Li W. CD-HIT: accelerated for clustering the next-generation sequencing data. Bioinformatics (Oxford, England). 2012;28(23):3150-2. Epub 2012/10/13. doi: 10.1093/bioinformatics/bts565. PubMed PMID: 23060610; PubMed Central PMCID: PMCPMC3516142.

39. Glez-Pena D, Gomez-Blanco D, Reboiro-Jato M, Fdez-Riverola F, Posada D. ALTER: program-oriented conversion of DNA and protein alignments. Nucleic Acids Res. 2010;38(Web Server issue):W14-8. Epub 2010/05/05. doi: 10.1093/nar/gkq321. PubMed PMID: 20439312; PubMed Central PMCID: PMCPMC2896128.

40. Suchard MA, Lemey P, Baele G, Ayres DL, Drummond AJ, Rambaut A. Bayesian phylogenetic and phylodynamic data integration using BEAST 1.10. Virus Evol. 2018;4(1):vey016. Epub 2018/06/27. doi: 10.1093/ve/vey016. PubMed PMID: 29942656; PubMed Central PMCID: PMCPMC6007674.

41. Rambaut A, Drummond AJ, Xie D, Baele G, Suchard MA. Posterior Summarization in Bayesian Phylogenetics Using Tracer 1.7. Syst Biol. 2018;67(5):901-4. Epub 2018/05/03. doi: 10.1093/sysbio/syy032. PubMed PMID: 29718447; PubMed Central PMCID: PMCPMC6101584.

42. Hu B, Zeng LP, Yang XL, Ge XY, Zhang W, Li B, et al. Discovery of a rich gene pool of bat SARS-related coronaviruses provides new insights into the origin of SARS coronavirus. PLoS Pathog. 2017;13(11):e1006698. Epub 2017/12/01. doi: 10.1371/journal.ppat.1006698. PubMed PMID: 29190287; PubMed Central PMCID: PMCPMC5708621.

43. Weiss SR, Leibowitz JL. Chapter 4 - Coronavirus Pathogenesis. In: Maramorosch K, Shatkin AJ, Murphy FA, editors. Advances in Virus Research. 81: Academic Press; 2011. p. 85-164.

44. Tong Y, Liu W, Liu P, Liu WJ, Wang Q, Gao GF. The origins of viruses: discovery takes time, international resources, and cooperation. Lancet. 2021;398(10309):1401-2. Epub 2021/10/04. doi: 10.1016/S0140-6736(21)02180-2. PubMed PMID: 34600605; PubMed Central PMCID: PMCPMC8483647 Correspondence.

45. Hu B, Ge X, Wang LF, Shi Z. Bat origin of human coronaviruses. Virol J. 2015;12:221. Epub 2015/12/23. doi: 10.1186/s12985-015-0422-1. PubMed PMID: 26689940; PubMed Central PMCID: PMCPMC4687304.

46. Li W, Shi Z, Yu M, Ren W, Smith C, Epstein JH, et al. Bats are natural reservoirs of SARS-like coronaviruses. Science. 2005;310(5748):676-9. Epub 2005/10/01. doi: 10.1126/science.1118391.

1          PubMed PMID: 16195424.

2   47.   Lau SK, Woo PC, Li KS, Huang Y, Tsoi HW, Wong BH, et al. Severe acute respiratory syndrome
3          coronavirus-like virus in Chinese horseshoe bats. Proc Natl Acad Sci U S A. 2005;102(39):14040-5.
4          Epub 2005/09/20. doi: 10.1073/pnas.0506735102. PubMed PMID: 16169905; PubMed Central
5          PMCID: PMCPMC1236580.

6   48.   Corman VM, Ithete NL, Richards LR, Schoeman MC, Preiser W, Drosten C, et al. Rooting the
7          phylogenetic tree of middle East respiratory syndrome coronavirus by characterization of a
8          conspecific virus from an African bat. J Virol. 2014;88(19):11297-303. Epub 2014/07/18. doi:
9          10.1128/JVI.01498-14. PubMed PMID: 25031349; PubMed Central PMCID: PMCPMC4178802.

10   49.   Wrobel AG, Benton DJ, Xu P, Roustan C, Martin SR, Rosenthal PB, et al. SARS-CoV-2 and bat RaTG13
11          spike glycoprotein structures inform on virus evolution and furin-cleavage effects. Nat Struct Mol
12          Biol. 2020;27(8):763-7. Epub 20200709. doi: 10.1038/s41594-020-0468-7. PubMed PMID:
13          32647346; PubMed Central PMCID: PMCPMC7610980.

14   50.   Murakami S, Kitamura T, Suzuki J, Sato R, Aoi T, Fujii M, et al. Detection and Characterization of
15          Bat Sarbecovirus Phylogenetically Related to SARS-CoV-2, Japan. Emerg Infect Dis.
16          2020;26(12):3025-9. Epub 2020/11/22. doi: 10.3201/eid2612.203386. PubMed PMID: 33219796;
17          PubMed Central PMCID: PMCPMC7706965.

18   51.   Delaune D, Hul V, Karlsson EA, Hassanin A, Ou TP, Baidaliuk A, et al. A novel SARS-CoV-2 related
19          coronavirus in bats from Cambodia. Nat Commun. 2021;12(1):6563. Epub 2021/11/11. doi:
20          10.1038/s41467-021-26809-4. PubMed PMID: 34753934; PubMed Central PMCID:
21          PMCPMC8578604.

22   52.   Zhou H, Ji J, Chen X, Bi Y, Li J, Wang Q, et al. Identification of novel bat coronaviruses sheds light
23          on the evolutionary origins of SARS-CoV-2 and related viruses. Cell. 2021;184(17):4380-91 e14.
24          Epub 2021/06/21. doi: 10.1016/j.cell.2021.06.008. PubMed PMID: 34147139; PubMed Central
25          PMCID: PMCPMC8188299.

26   53.   Liu K, Pan X, Li L, Yu F, Zheng A, Du P, et al. Binding and molecular basis of the bat coronavirus
27          RaTG13 virus to ACE2 in humans and other species. Cell. 2021;184(13):3438-51 e10. Epub
28          20210524. doi: 10.1016/j.cell.2021.05.031. PubMed PMID: 34139177; PubMed Central PMCID:
29          PMCPMC8142884.

30   54.   Liu K, Tan S, Niu S, Wang J, Wu L, Sun H, et al. Cross-species recognition of SARS-CoV-2 to bat ACE2.
31          Proc Natl Acad Sci U S A. 2021;118(1):DOI: 10.1073/pnas.2020216118. Epub 2020/12/19. doi:
32          10.1073/pnas.2020216118. PubMed PMID: 33335073; PubMed Central PMCID: PMCPMC7817217.

33   55.   Gao GF, Wang L. COVID-19 Expands Its Territories from Humans to Animals. China CDC Wkly.
34          2021;3(41):855-8. Epub 2021/10/28. doi: 10.46234/ccdcw2021.210. PubMed PMID: 34703641;
35          PubMed Central PMCID: PMCPMC8521158.

36   56.   Liu P, Jiang JZ, Wan XF, Hua Y, Li L, Zhou J, et al. Are pangolins the intermediate host of the 2019
37          novel coronavirus (SARS-CoV-2)? PLoS Pathog. 2020;16(5):e1008421. Epub 20200514. doi:
38          10.1371/journal.ppat.1008421. PubMed PMID: 32407364; PubMed Central PMCID:

1       PMCPMC7224457.

2    57.  Jones DT, Taylor WR, Thornton JM. The rapid generation of mutation data matrices from protein

3        sequences. Comput Appl Biosci. 1992;8(3):275-82. Epub 1992/06/01. doi:

4        10.1093/bioinformatics/8.3.275. PubMed PMID: 1633570.

5    58.  Tamura K, Nei M. Estimation of the number of nucleotide substitutions in the control region of

6        mitochondrial DNA in humans and chimpanzees. Mol Biol Evol. 1993;10(3):512-26. Epub

7        1993/05/01. doi: 10.1093/oxfordjournals.molbev.a040023. PubMed PMID: 8336541.

8    59.  Saitou N, Nei M. The neighbor-joining method: a new method for reconstructing phylogenetic

9        trees. Mol Biol Evol. 1987;4(4):406-25. Epub 1987/07/01. doi:

10        10.1093/oxfordjournals.molbev.a040454. PubMed PMID: 3447015.

11    60.  Tamura K, Nei M, Kumar S. Prospects for inferring very large phylogenies by using the neighbor-

12        joining method. Proceedings of the National Academy of Sciences (USA). 2004;101:11030-5.

13

14

15  **Figure legend**

16  **Figure 1.** Accumulation of sequenced genomes, confirmed cases and mutations of SARS-CoV-2.

17  A, Accumulation of sequenced genomes, confirmed cases. B, Correlation of sequenced genomes

18  and confirmed cases. C, Accumulation curve of nucleotide diversity (Pi) in genomes as calculated

19  using pairwise deletion (PD) model; D, Accumulation curve of theta(W).

20  **Figure 2. Mutation accumulation and phylogenetic analysis of NS8 variants in SARS-CoV-2**.

21  A and B, Genome accumulation curve of NS8_84L and NS8_84S as counted by the sample

22  collection date; C and D, Correlation analysis between genomes of genotypes and confirmed cases;

23  E and F, Mutation (Pi) accumulation curves of genotype NS8_84L (E) and NS8_84S (F); G,

24  Maximum Likelihood tree of peptide sequences of NS8 as inferred by using the JTT matrix-based

25  model [57]. The bootstrap supports by 1000 replicates are shown next to the branches. H, Maximum

26  Likelihood tree of gene sequences of NS8 as inferred by using the Tamura-Nei model [58]. The

27  bootstrap supports are shown above the branches. All positions with less than 95% site coverage

28  were eliminated. Evolutionary analyses were conducted in MEGA11 [36].

29  **Figure 3. Mutation accumulation and phylogenetic analysis of Spike variants in SARS-CoV-**

30  **2**. A and B, Genome accumulation curve of S_614G and S_614D as counted by the sample

31  collection date; C and D, Correlation analysis between genomes of S_614G and S_614D and

22

1   confirmed cases; E and F, Mutation (Pi) accumulation curves of genotype S_614G (E) and S_614D

2   (F); G, Maximum Likelihood tree of gene sequences of Spike as inferred by using the Tamura-Nei

3   model [58]; H, Evolutionary tree of peptide sequences of Spike as inferred by using the Neighbor-

4   Joining method [59]. The bootstrap supports by 1000 replicates are shown next to the branches. All

5   positions with less than 95% site coverage were eliminated. Evolutionary analyses were conducted

6   in MEGA11 [36].

7   **Figure 4. Analysis of main haplotypes of SARS-CoV-2**. A-C, Genome accumulation curves of

8   haplotype GL (A), GS (B), DL (C), and DS (C); D, Haplotype profile in China, Europe, USA, and

9   other parts of the world; E-H, Correlation analysis between genomes of main haplotypes and

10  confirmed cases; I-L, Mutation (Pi) accumulation curves of main genotypes.

11  **Figure 5. Phylogenetic analysis of main haplotypes of SARS-CoV-2**. The trees were inferred by

12  using Maximum Likelihood (A) and Neighbor-Joining (B) methods and rooted with the bat SARS-

13  related viral genome RaTG13. For ML tree, Tamura-Nei model [58] was used, all positions with

14  less than 95% site coverage were eliminated, and the tree with the highest log likelihood (-48997.72)

15  is shown. For NJ tree, the evolutionary distances were computed using the Maximum Composite

16  Likelihood method [60], and all ambiguous positions were removed for each sequence pair

17  (pairwise deletion option). The bootstrap supports by 1000 replicates are shown next to the branches.

18  All positions with less than 95% site coverage were eliminated. The percentages of replicate trees

19  in which the associated taxa clustered together in the bootstrap test (1000 replicates) are shown next

20  to the branches.  Evolutionary analyses were conducted in MEGA11 [36].

21  **Figure 6.** Posterior distribution for the tMRCA of haplotypes DL, DS, GL, and their combination

22  (DS+DL+GL). Redundant genomes were removed by using CD-HIT {Fu, 2012 #1002} with a

23  threshhold of 99.97% identities. Inference was performed by using a strict molecular clock and a

24  Bayesian Skyline coalescent prior in BEAST 1.10 [40], and summarized with Tracer 1.7 [41]. The

25  mean tMRCA is shown above or beside the curves.

26  **Figure 7. Evolutionary trajectory of SARS-CoV-2 haplotypes.** The codons of S_614 and NS8_84 are

27  provided below the haplotypes. The single nucleotides mutated from previous haplotypes were

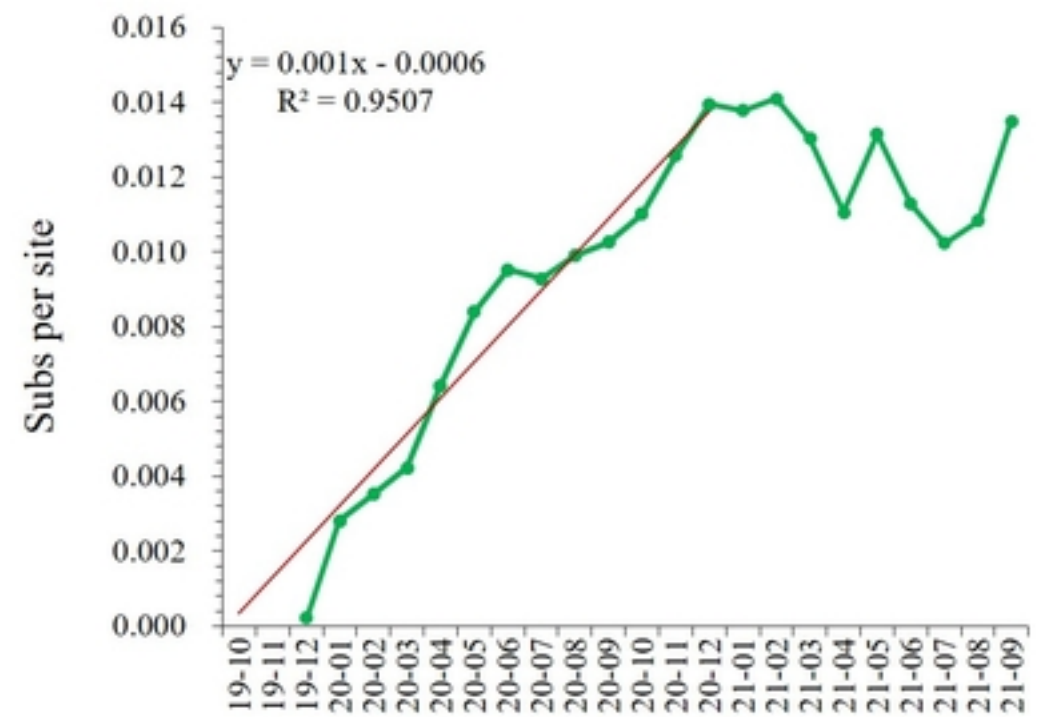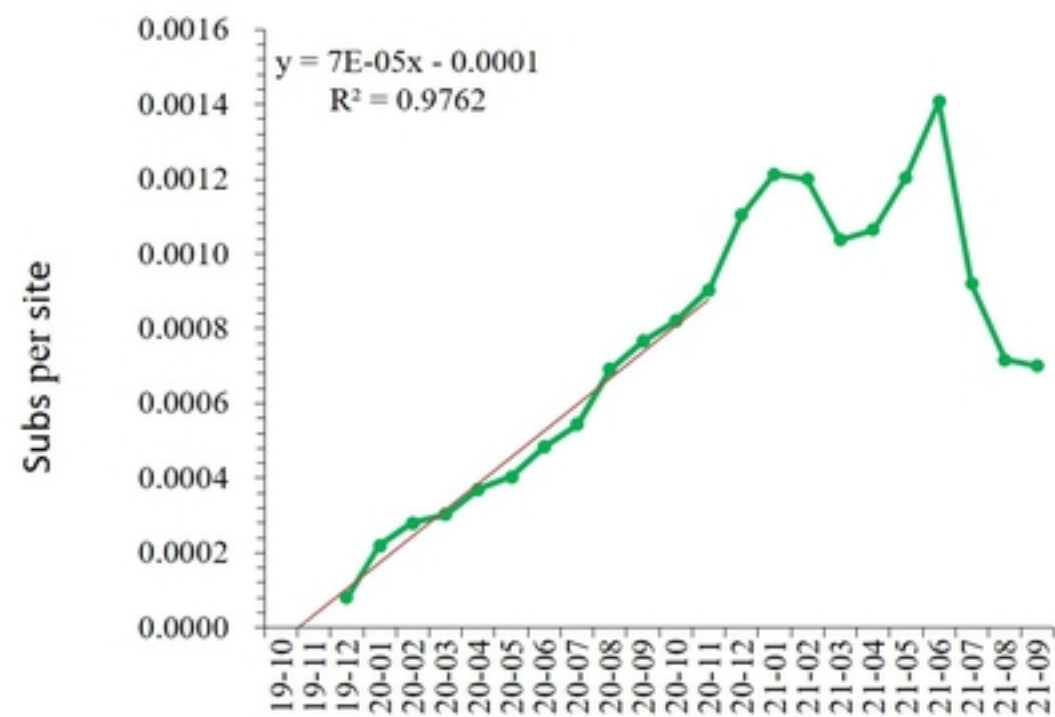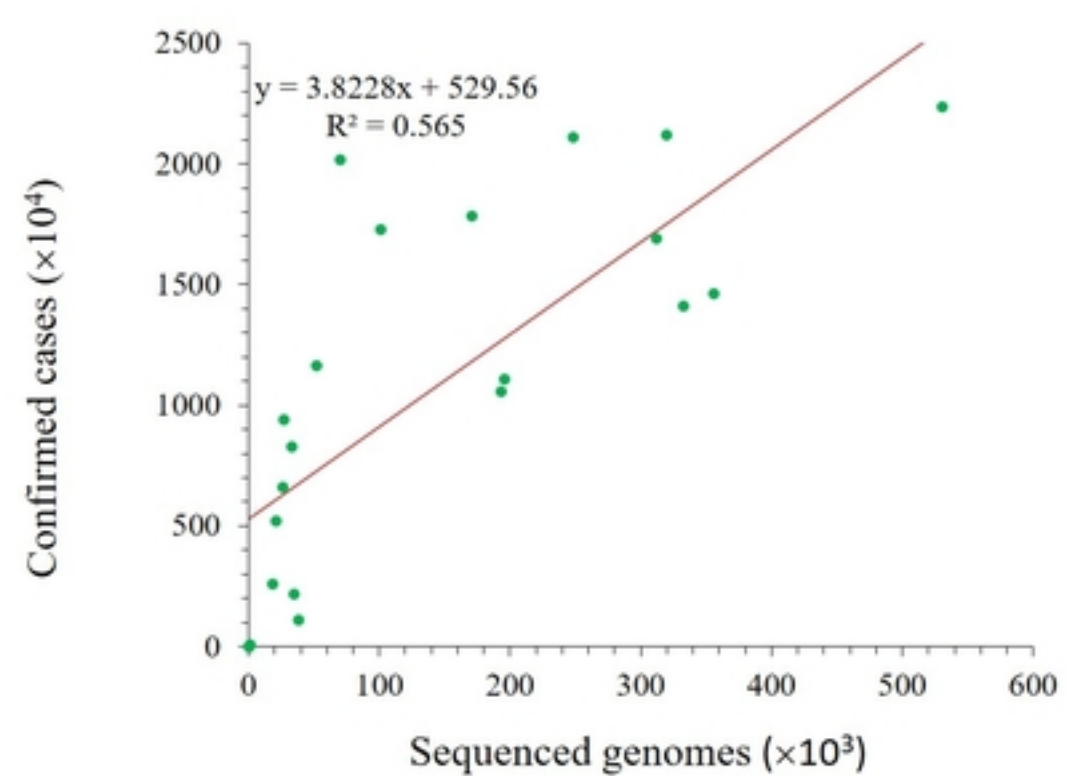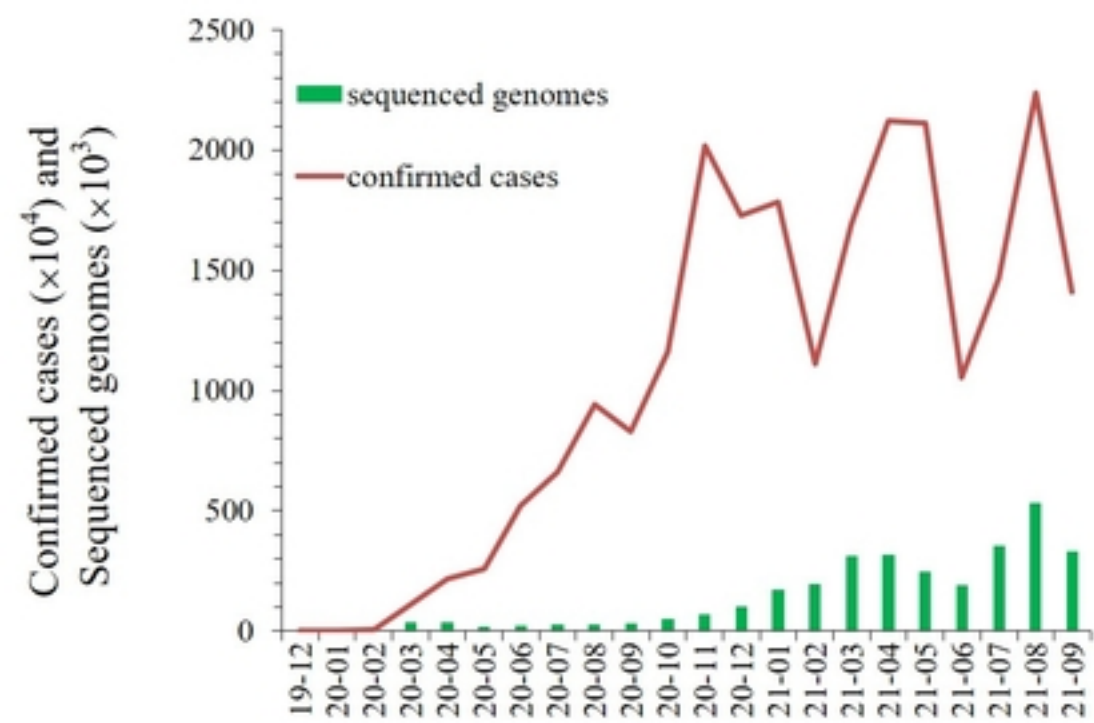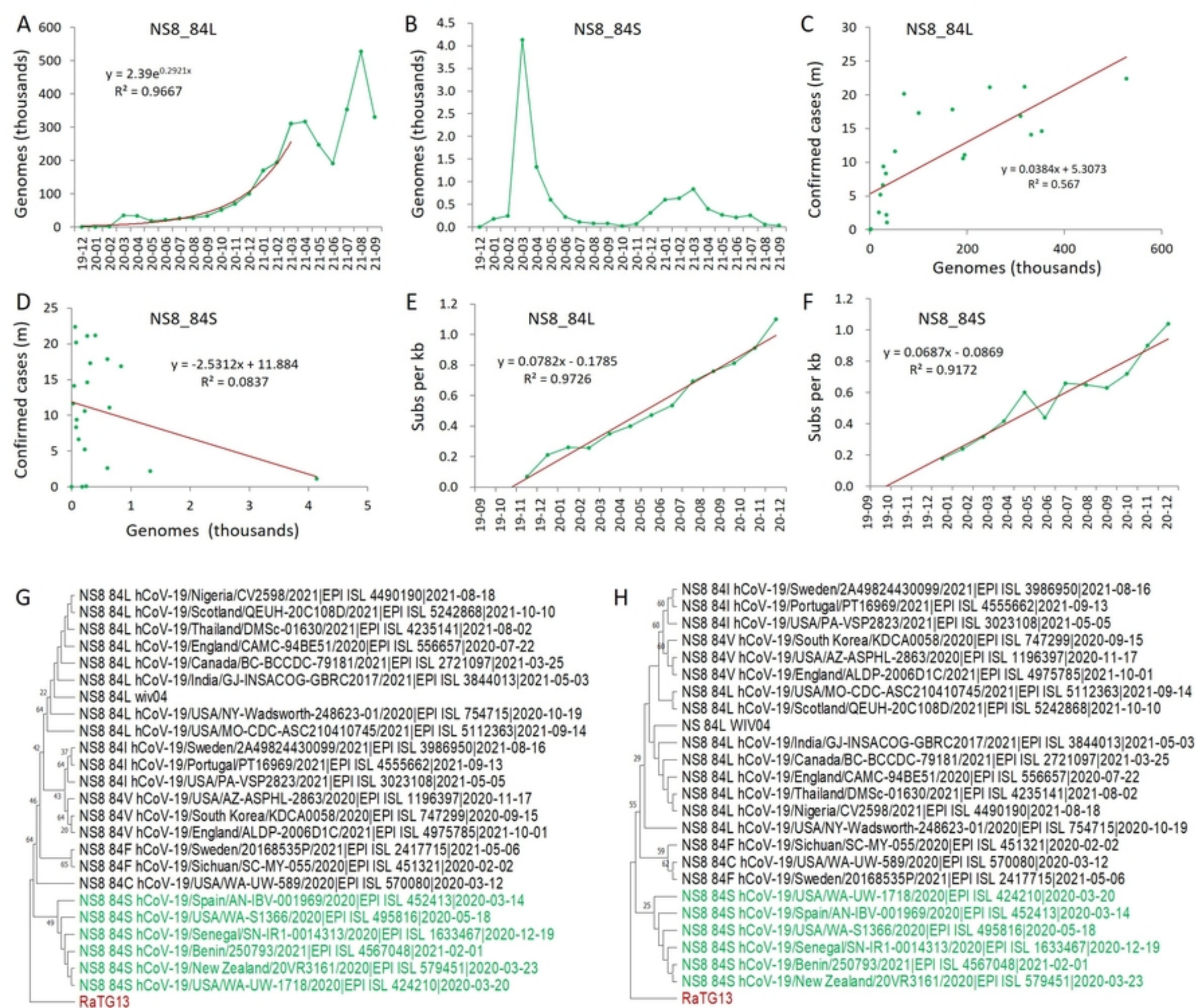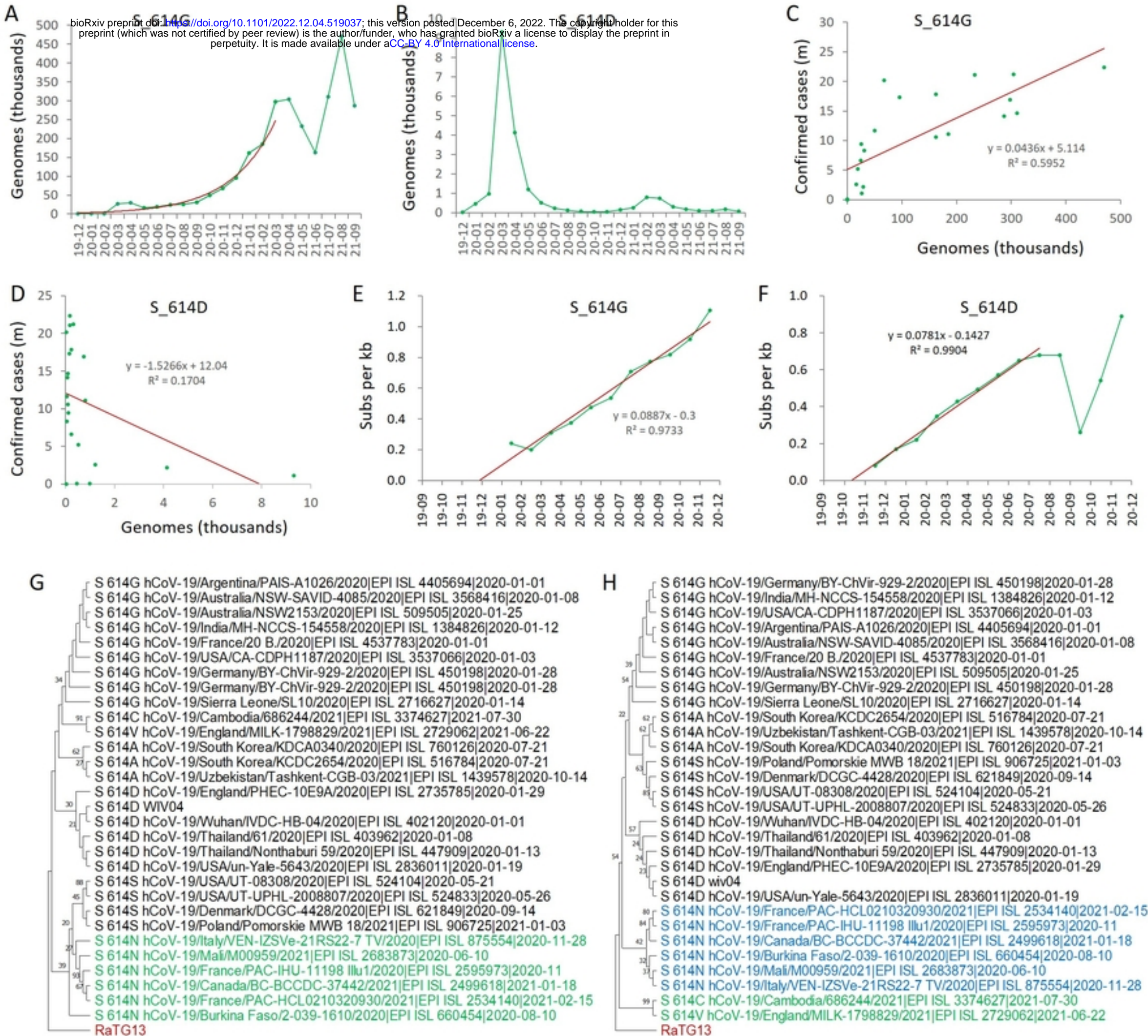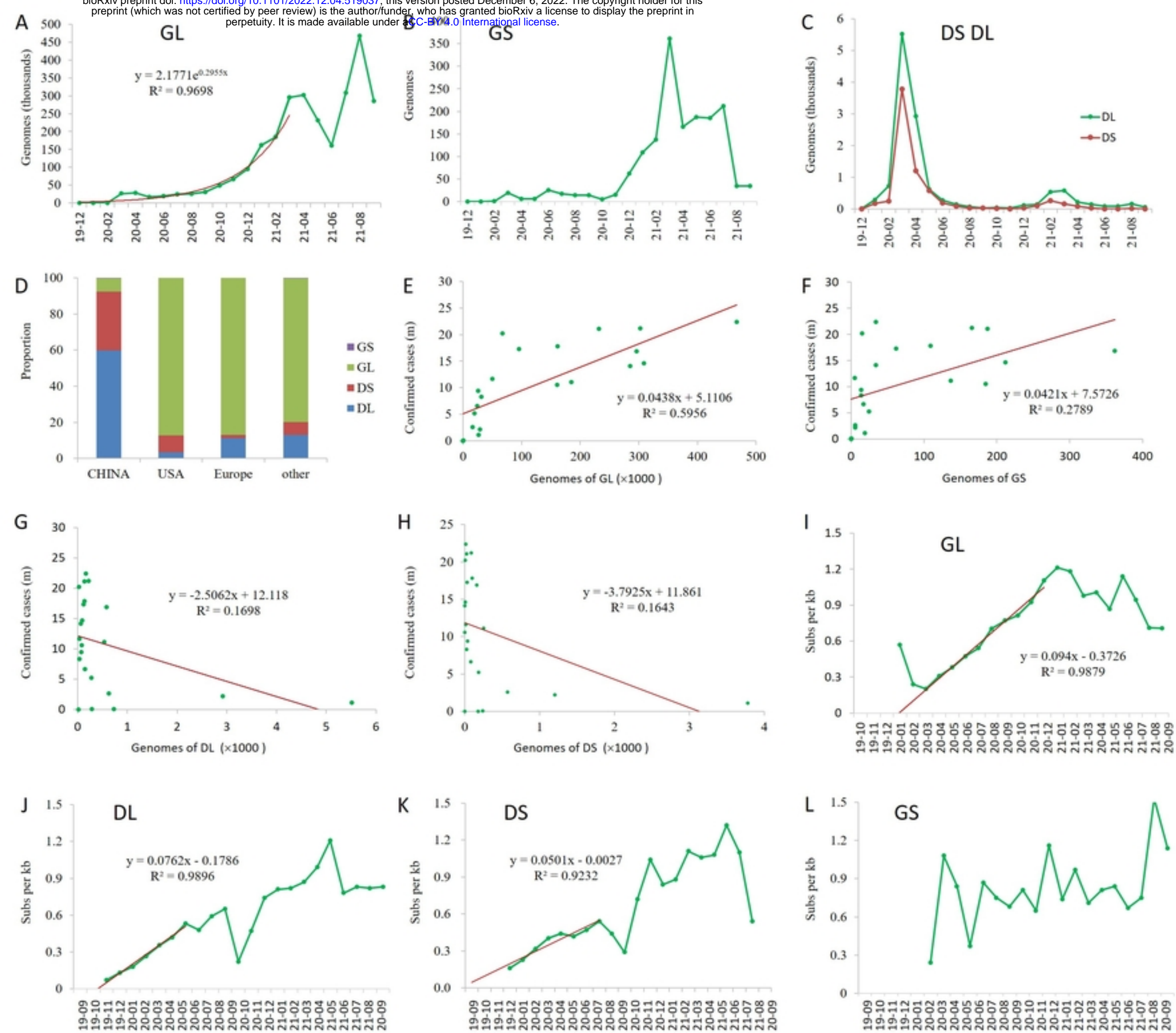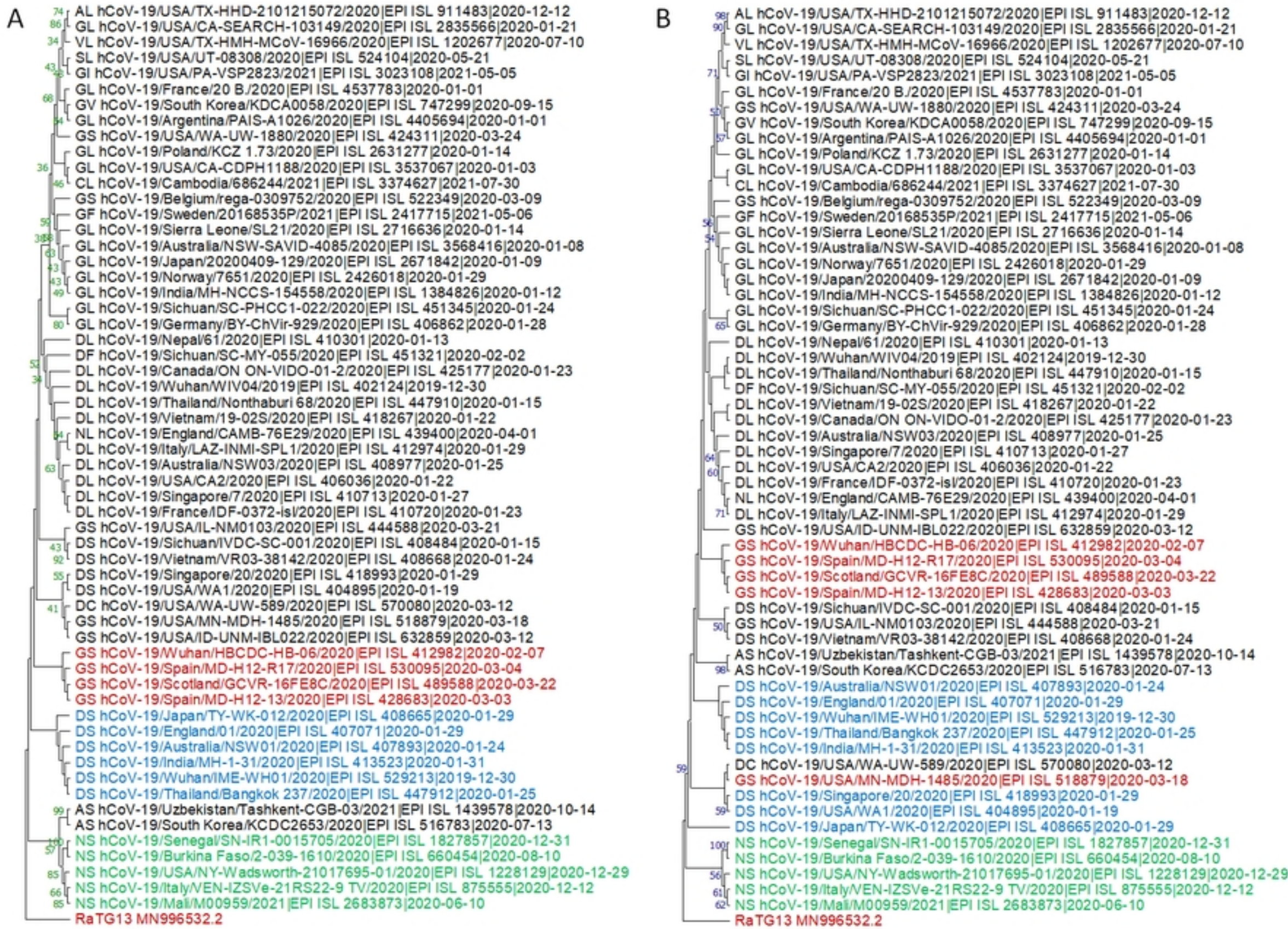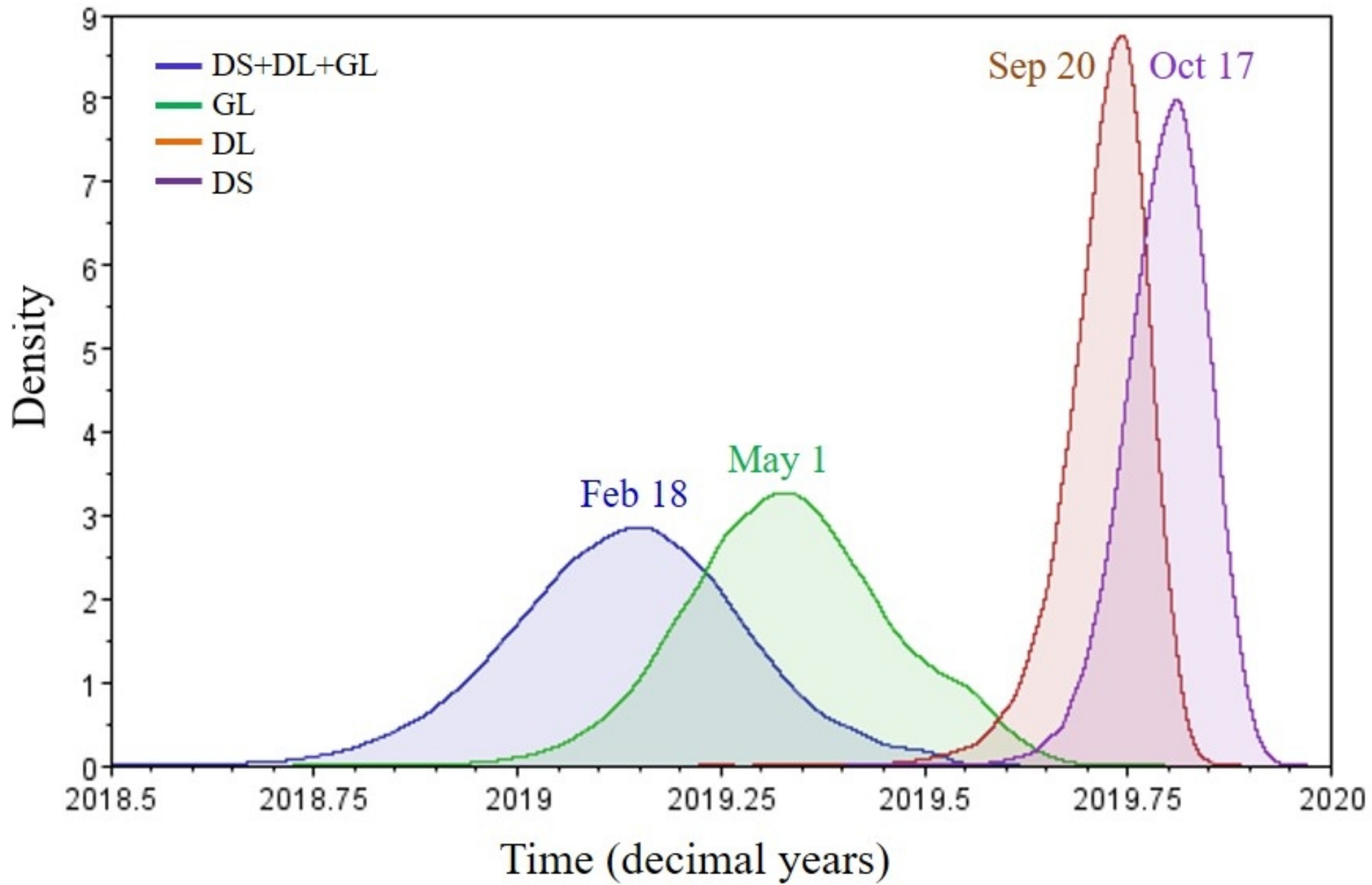28  high-lighted in green. The proposed time axis is provided below.
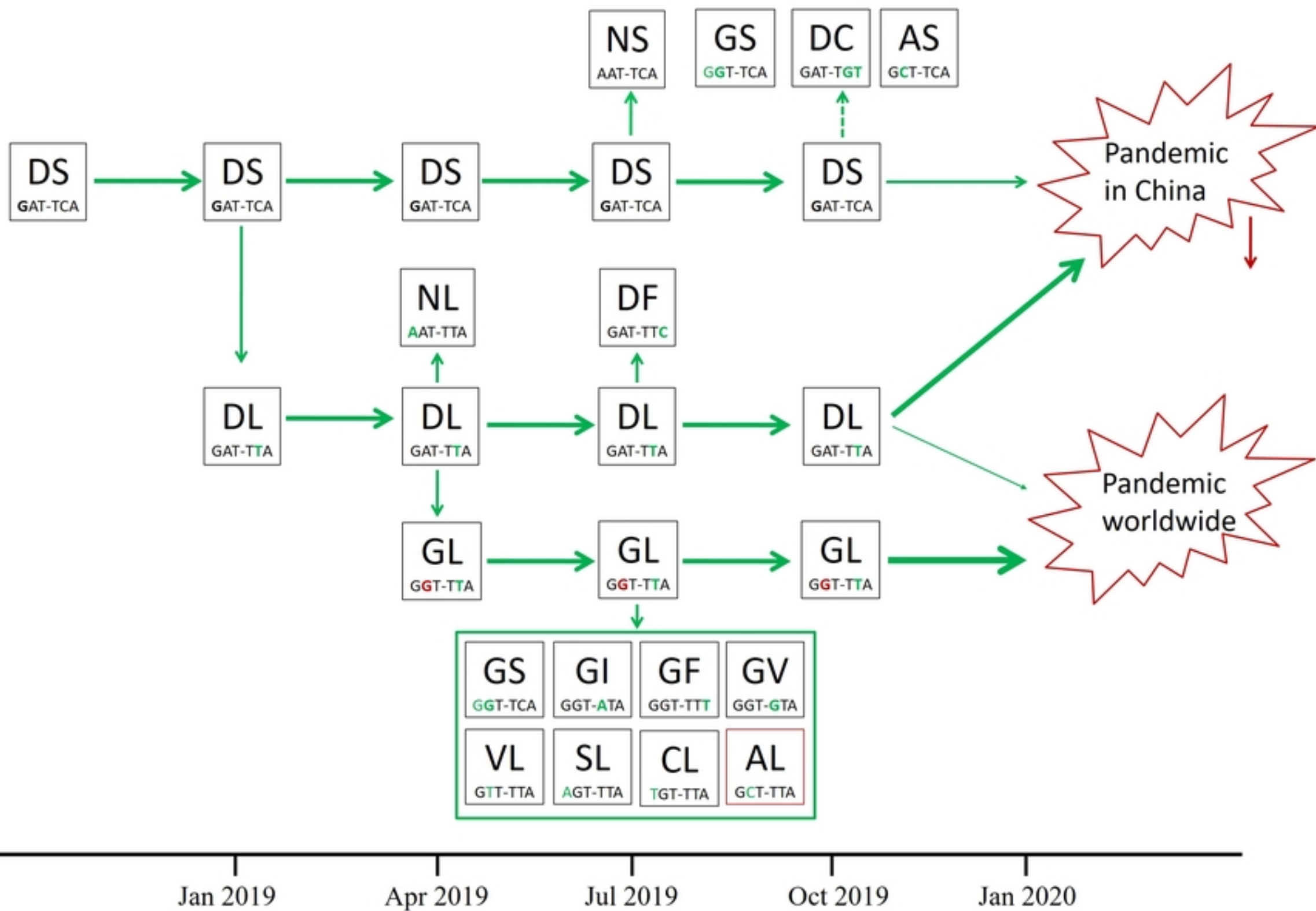
29

Figure 1

Figure 2

Figure 3

Figure 4

Figure 5

Figure 6

Figure 7