

Synonymous mutations and the molecular evolution of SARS-Cov-2 origins

Hongru Wang¹, Lenore Pipes¹ and Rasmus Nielsen^{1, 2, 3*}

¹ Department of Integrative Biology, UC Berkeley, Berkeley, CA 94707, USA.

²Department of Statistics, UC Berkeley, Berkeley, CA 94707, USA.

³Globe Institute, University of Copenhagen, 1350 København K, Denmark.

*Address: 4098 Valley Life Sciences Building, Department of Integrative Biology, UC Berkeley. Berkeley, CA 94707. rasmus_nielsen@berkeley.edu

1 **Abstract**

2 Human severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) is most closely
3 related, by average genetic distance, to two coronaviruses isolated from bats, RaTG13 and
4 RmYN02. However, there is a segment of high amino acid similarity between human SARS-
5 CoV-2 and a pangolin isolated strain, GD410721, in the receptor binding domain (RBD) of
6 the spike protein, a pattern that can be caused by either recombination or by convergent
7 amino acid evolution driven by natural selection. We perform a detailed analysis of the
8 synonymous divergence, which is less likely to be affected by selection than amino acid
9 divergence, between human SARS-CoV-2 and related strains. We show that the
10 synonymous divergence between the bat derived viruses and SARS-CoV-2 is larger than
11 between GD410721 and SARS-CoV-2 in the RBD, providing strong additional support for
12 the recombination hypothesis. However, the synonymous divergence between pangolin
13 strain and SARS-CoV-2 is also relatively high, which is not consistent with a recent
14 recombination between them, instead it suggests a recombination into RaTG13. We also
15 find a 14-fold increase in the d_N/d_S ratio from the lineage leading to SARS-CoV-2 to the
16 strains of the current pandemic, suggesting that the vast majority of non-synonymous
17 mutations currently segregating within the human strains have a negative impact on viral
18 fitness. Finally, we estimate that the time to the most recent common ancestor of SARS-
19 CoV-2 and RaTG13 or RmYN02 based on synonymous divergence, is 51.71 years (95%
20 C.I., 28.11-75.31) and 37.02 years (95% C.I., 18.19-55.85), respectively.

21

22 **Introduction**

23 The Covid19 pandemic is perhaps the biggest public health and economic threat that the world
24 has faced for decades (Li, et al. 2020; Wu, et al. 2020; Zhou, Yang, et al. 2020). It is caused by

1 a coronavirus (Lu, et al. 2020; Zhang and Holmes 2020), Severe acute respiratory syndrome
2 coronavirus 2 (SARS-CoV-2), an RNA virus with a 29,891 bp genome consisting of four major
3 structural genes (Wu, et al. 2020; Zhou, Yang, et al. 2020). Of particular relevance to this study
4 is the *spike* protein which is responsible for binding to the primary receptor for the virus,
5 angiotensin-converting enzyme 2 (*ACE2*) (Wan, et al. 2020; Wu, et al. 2020; Zhou, Yang, et al.
6 2020).

7 Human SARS-CoV-2 is related to a coronavirus (RaTG13) isolated from the bat
8 *Rhinolophus affinis* from Yunnan province of China (Zhou, Yang, et al. 2020). RaTG13 and the
9 human strain reference sequence (Genbank accession number MN996532) are 96.2% identical
10 and it was first argued that, throughout the genome, RaTG13 is the closest relative to human
11 SARS-CoV-2 (Zhou, et al. 2020). Zhang, et al. 2020 showed that RaTG13 and SARS-CoV-2
12 were 91.02% and 90.55% identical ,respectively, to coronaviruses isolated from pangolins
13 (Pangolin-CoV), which therefore form a close outgroup to the SARS-CoV-2+RaTG13 clade .
14 Furthermore, five key amino acids in the receptor-binding domain (RBD) of *spike* were identical
15 between SARS-CoV-2 and Pangolin-CoV, but differed between those two strains and RaTG13.
16 (Lam, et al. 2020) independently made similar observations and additionally showed that when
17 analyzing a window of length 582bp in the RBD, nonsynonymous mutations support a
18 phylogenetic tree with SARS-CoV-2 and Pangolin-CoV as sister-groups, while synonymous
19 mutations do not. They discuss two possible explanations for their results, one which includes
20 recombination and another which includes selection-driven convergent evolution. Boni, et al.
21 2020 found little evidence of recombination between SARS-CoV-2 and Pangolin-CoV using the
22 recombination analysis software 3SEQ (Lam, et al. 2018), but argued that if there has been
23 recombination, it likely occurred into RaTG13 from an unknown divergent source. This would
24 explain the amino acid similarity between SARS-CoV-2 and Pangolin-CoV in the RBD as an
25 ancestral trait that has been lost (by recombination) in RaTG13. Using a phylogenetic analysis
26 they also dated the RaTG13 and SARS-CoV-2 divergence to be between 40 to 70 years.

1 Recently, Zhou et al. (2020) discovered a viral strain (Zhou, et al. 2020), RmYN02 from the bat
2 *Rhinolophus malayanus*, with a reported 97.2% identity in the ORF1ab gene but with only
3 61.3% sequence similarity to SARS-CoV-2 in the RBD. Moreover, the RmYN02 strain also
4 harbors multiple amino acid insertions at the functional polybasic (furin) cleavage site (Zhou, et
5 al. 2020) that was thought to be unique to the human SARS-CoV-2 (Andersen, et al. 2020).

6 To analyze the history of these sequences further, we here focus on patterns of
7 synonymous divergence, which has received less focus, but also is less likely to be affected by
8 selection than amino acid divergence. We develop a bias corrected estimator of synonymous
9 divergence specific for SARS-CoV-2 and related strains, and analyze divergence using both
10 sliding windows and a whole-genome approach between SARS-CoV-2 and related viral strains.

11

12 **Materials and methods**

13 *Sequence data:* The pangolin virus sequences, GD410721 and GX_P1E, were downloaded
14 from GISAID with accession numbers EPI_ISL_410721 and EPI_ISL_410539, respectively, and
15 RmYN02 was provided by E. C. Holmes. All other sequences analyzed in this study were
16 downloaded from either NCBI Genbank or National Microbiology Data Cente (NMDC). The
17 accession codes for non-human sequences can be found in Supplementary Table 2 and the
18 accession codes for human sequences can be found in Supplementary Table 3.

19

20 *BLAST searches:* Sequences for blast databases were downloaded on March 26, 2020 from the
21 following sources: EMBL nucleotide libraries for virus

22 (<ftp://ftp.ebi.ac.uk/pub/databases/embl/release/std>), NCBI Virus Genomes

23 (<ftp://ftp.ncbi.nlm.nih.gov/genomes/Viruses>), NCBI Virus Genbank Entries

24 (<ftp://ftp.ncbi.nlm.nih.gov/genomes/genbank/viral/>), NCBI Influenza Genomes

25 (<ftp://ftp.ncbi.nlm.nih.gov/genomes/INFLUENZA/>), all Whole Genome Shotgun

26 (<https://www.ncbi.nlm.nih.gov/genbank/wgs/>) assemblies under taxonomy ID 10239, along with

1 GISAID Epiflu and EpiCoV databases. Sequences from Valitutto et al. (2020) were also added
2 to the database. Blast databases were created using the default parameters for makeblastdb.
3 Blast searches were performed using blastn with parameters “-word_size 7 -reward 1 -penalty -
4 3” and all other parameters as the default settings. All the blast hits to different Guangdong
5 pangolin viral strain sequences were merged as one hit, and the blast hits to different Guangxi
6 pangolin viral strain sequences were also merged.

7
8 *Alignment:* To obtain an in-frame alignment of the genomes, we first identified the coding
9 sequences of each viral strain using independent pairwise alignments with the coding
10 sequences of the SARS-CoV-2 (Wuhan-Hu-1) genome. The alignments were performed using
11 MAFFT (v7.450) (Kato and Standley 2013) with parameters “--maxiterate 1000 --localpair”. We
12 then performed multiple sequence alignments for the coding sequences of each gene using
13 PRANK (Loytynoja 2014) (v.170427) with parameters “-codon -F”. Finally, the alignments for all
14 genes were concatenated following their genomic order. ORF1a was excluded without loss of
15 any nucleotides, since its sequence is a subset of ORF1ab.

16
17 *Recombination detection:* We detected possible recombination events across the genome using
18 the 3SEQ algorithm (Lam, et al. 2018) implemented in RDP (Martin, et al. 2015) (version Beta
19 5.5). The analysis was performed on the multiple sequence alignment consisting of the five viral
20 strains with the 'automated 3Seq' mode in RDP. All regions in which any of the viral strains
21 showed evidence ($p < 0.05$) of recombination (Supplementary Table 5) were removed in
22 subsequent analyses from all strains when stating that recombination regions were removed.

23
24 *Tree estimation:* We estimated phylogenetic trees using two methods: Neighbor Joining (NJ)
25 and Maximum Likelihood (ML). The NJ trees were estimated using d_N or d_S distance matrices
26 which estimated using codeml (Yang 2007) with parameters " runmode= -2, CodonFreq = 2,

1 cleandata = 1". To obtain bootstrap values, we bootstrapped the multiple sequence alignments
2 1,000 times, repeating the inference procedure for each bootstrap sample. The NJ tree was
3 estimated using the 'neighbor' software from the PHYLIP package (Felsenstein 2009). For ML
4 trees, we used IQ-TREE (Nguyen, et al. 2015) (v1.5.2) with parameter "-TEST -alrt 1000" which
5 did substitution model selection for the alignments and performed maximum-likelihood tree
6 estimation with the selected substitution model for 1,000 bootstrap replicates. For this analysis,
7 we masked all regions (Supplementary Table 5) identified by 3SEQ (Lam, et al. 2018) to be
8 recombinant in any of the five studied viral genome. The coordinates (based on the Wuhan-Hu-
9 1 genome) of the three recombination regions (merged set of all the regions in Supplementary
10 Table 5) were: 14611-15225, 21225-24252 and 25965-28297. We also estimate genome-wide
11 divergence between RaTG13 and Wuhan-Hu-1 only excluding the region (position 22853-23092)
12 where potential recombination was detected for the Wuhan-Hu-1 strain (Supplementary Table
13 5).

14

15 *Simulations:* We simulated divergence with realistic parameters for SARS-CoV-2 using a
16 continuous time Markov chain under the F3x4 codon-based model (Goldman and Yang 1994;
17 Muse and Gaut 1994; Yang, et al. 2000), which predicts codon frequencies from the empirical
18 nucleotide frequencies in all 3 codon positions and using the global genomic maximum
19 likelihood estimates of the transition/transversion bias κ (=2.9024) and the d_N/d_S ratio ω
20 (=0.0392) estimated from the human SARS-CoV-2 comparison to the nearest outgroup
21 sequence, RaTG13 (see Results). For the simulations of short 300 bp sequences we kept ω
22 constant but varied time such that the number of synonymous substitutions per synonymous
23 sites, d_S , varied between 0.25 and 3.00. Estimates of $d_S > 3$ are truncated to 3. For simulations
24 of genome-wide divergence between RaTG13 and human strains, we fix d_S at 0.1609 (the
25 maximum likelihood estimate outside the RBD region reported in the Results section). In all
26 cases, we use 10,000 independent replicate simulations for each parameter setting.

1
2 *Estimation of sequence divergence in 300-bp windows:* d_N and d_S were estimated using two
3 different methods implemented in the PAML package (Yang 2007) (version 4.9d): a count-
4 based method, YN00 (Yang and Nielsen 2000) as implemented in the program 'yn00' with
5 parameters "icode = 0, weighting = 0, commonf3x4 = 0", and a maximum-likelihood method
6 (Goldman and Yang 1994; Muse and Gaut 1994) implemented in codeml applied with
7 arguments "runmode= -2, CodonFreq = 2". The estimates in 300-bp windows were further bias-
8 corrected as described below.

9
10 *Bias correction for d_S estimates in 300-bp window:* To correct for the biases observed in the
11 estimation of d_S (see results section) we identified a quartic function which maps from \hat{d}_S , the
12 estimates of d_S , into \widehat{d}_S^* , the bias corrected estimate such that to a close approximation, $E[\widehat{d}_S^*]$
13 = d_S . To identify the coefficients of this function we used 10,000 simulations as previously
14 described, on a grid of d_S values (0.25, 0.5, 0.75, ..., 3.0). We then identified coefficients such
15 that sum of $(E[\widehat{d}_S^*] - d_S)^2$ is minimized over all simulation values.

16

17 **Results**

18 *Database searches*

19 To identify possible viral strains that may have contributed, by recombination, to the formation of
20 human SARS-CoV-2, we searched NCBI and EMBL virus entries along with GISAID Epiflu and
21 EpiCov databases for similar sequences using BLAST in 100bp windows stepping every 10bp
22 (Fig. 1B). The majority of the genome (78.1%, 2330/2982 of the windows) has one unique best
23 hit, likely reflecting the high genetic diversity of the coronavirus. 21.9% of the genomic regions
24 has multiple best hits, which suggests that these regions might be more conserved. Among the
25 windows with unique best hits, 97.0% (2260/2330) of them were the RaTG13 or RmYN02 bat

1 strains and 1.9% of them, including the *ACE2* contact residues region of the S protein, were the
2 pangolin SARS-CoV-2 virus. These observations are consistent with previous results that
3 RaTG13 and RmYN02 are the most closely related viral strains, while the region containing the
4 *ACE2* contact residues is more closely related to the pangolin virus strains. A considerable
5 amount of genomic regions (20 windows with unique hits) show highest sequence identity with
6 other coronaviruses of the SARS-CoV-2 related lineage (Lam, et al. 2020) (bat-SL-CoVZC45
7 and bat-SL-CoVZXC21, (Hu, et al. 2018)). In addition, there were 6 windows whose unique top
8 hits are coronavirus of a SARS-CoV related lineage (Lam, et al. 2020) (Supplementary Table 4).
9 The mosaic pattern that different regions of the genome show highest identity to different virus
10 strains is likely to have been caused by the rich recombination history of the SARS-CoV-2
11 lineage (Boni, et al. 2020; Patiño-Galindo, et al. 2020). Moreover, its unique connection with
12 SARS-CoV related lineages in some genomic regions may suggest recombination between the
13 ancestral lineage of SARS-CoV-2 and distantly related virus lineages, although more formal
14 analyses are needed to determine the recombination history (see also Boni, et al. 2020 for
15 further discussion). Searching databases with BLAST using the most closely related viral
16 strains, RaTG13 and RmYN02, we observe a very similar pattern, as that observed for SARS-
17 CoV-2, in terms of top hits across the genome (Fig. 1B), suggesting that these possible
18 recombination events with distantly related lineages are not unique to the SARS-CoV-2 lineage,
19 but happened on the ancestral lineage of both SARS-CoV-2, RaTG13, and RmYN02. A notable
20 exception is a large region around the S gene, where RmYN02 show little similarity to both
21 SARS-CoV-2 and RaTG13.

22

23 *Sequence similarity*

24 We focus further on analyzing SARS-CoV-2: Wuhan-Hu-1 as the human nCoV19 reference
25 strain (Wu, et al. 2020) and the four viral strains with highest overall identity: the bat strains ,
26 RmYN02 and RaTG13 (Zhou, et al. 2020; Zhou, Yang, et al. 2020) and the Malayan pangolin

1 strains, GD410721 and GX_P1E, which were isolated from samples seized by Guangdong and
2 Guangxi Customs of China, respectively (Lam, et al. 2020; Xiao, et al. 2020). These four strains
3 have previously been identified as the strains most closely related to SARS-CoV-2 (Lam, et al.
4 2020; Zhou, et al. 2020). The overall Neighbor-Joining and maximum-likelihood trees between
5 the sequences closely related to human SARS-CoV-2 is shown in Fig. 2. RmYN02 has
6 previously been reported to have discordant levels of divergence to SARS-CoV-2 in different
7 parts of the genome, presumably as a result of recombination. To illustrate this, we performed
8 recombination analyses with 3SEQ (Lam, et al. 2018) across the five viral genomes and
9 identified several recombination regions affecting RmYN02 (Fig. 1C; Supplementary Table 5).
10 Phylogenetic analysis (Fig. 2) in genome region with all recombination tracts (Supplementary
11 Table 5) masked using Maximum-likelihood tree (Fig. 2A) and Neighbor-joining tree based on
12 synonymous (Fig. 2B) or non-synonymous (Fig. 2C) mutation distance metrics consistently
13 support RmYN02 as the nearest outgroup to human SARS-CoV-2, in contrast to previous
14 analyses before the discovery of RmYN02, which instead found RaTG13 to be the nearest
15 outgroup (Lam, et al. 2020; Zhou, Yang, et al. 2020). This observation is also consistent with the
16 genome-wide phylogeny constructed in Zhou et al (Zhou, et al. 2020).

17 We plot the overall sequence similarity (% nucleotides identical) between SARS-CoV-2
18 and the four other strains analyzed in windows of 300 bp (Fig. 1). Notice that the divergence
19 between human SARS-CoV-2 and the bat viral sequences, RaTG13 and RmYN02, in most
20 regions of the genome, is quite low compared to the other comparisons. A notable exception is
21 the suspected recombination region in RmYN02 that has an unusual high level of divergence
22 with all other viruses (Fig. 2E). However, there is also another exception: a narrow window in
23 the RBD of the S gene where the divergence between SARS-CoV-2 and GD410721 is
24 moderate and the divergences between GD410721 and both SARS-CoV-2 and RaTG13 are
25 quite high and show very similar pattern. This would suggest a recombination event from a
26 strain related to GD410721 into an ancestor of the human strain (Lam, et al. 2020; Xiao, et al.

1 2020; Zhang, et al. 2020), or alternatively, from some other species into RaTG13, as previously
2 hypothesized (Boni, et al. 2020). We note that RmYN02 is not informative about the nature of
3 this event as it harbors a long and divergent haplotype in this region, possibly associated with
4 another independent recombination event with more distantly related viral strains (Fig. 2E). The
5 other four sequences are all highly, and approximately equally, divergent from RmYN02 in this
6 large region (Fig. 2E), suggesting that the RmYN02 strain obtained a divergent haplotype from
7 the recombination event. When BLAST searching using 100-bp windows along the RmYN02
8 genome, we find no single viral genome as the top hit, instead the top hits are found
9 sporadically in different viral strains of the SARS-CoV lineage (Fig. 2F), suggesting that the
10 sequence of the most proximal donor is not represented in the database.

11

12 *Estimating synonymous divergence and bias correction*

13 While the overall divergence in the *S* gene encoding the *spike* protein could suggest the
14 presence of recombination in the region, Lam *et al.* (2020) (Lam, et al. 2020) reported that the
15 tree based on synonymous substitutions supported RaTG13 as the sister taxon to the human
16 SARS-CoV-2 also in this region. That would suggest the similarity between GD410721 and
17 human SARS-CoV-2 might be a consequence of convergent evolution, possibly because both
18 strains adapted to the use of the same receptor. An objective of the current study is to examine
19 if there are more narrow regions of the spike protein that might show evidence of recombination.
20 We investigate this issue using estimates of synonymous divergence per synonymous site (d_S)
21 in sliding windows of 300 bp. However, estimation of d_S is complicated by the high levels of
22 divergence and extremely skewed nucleotide content in the 3rd position of the sequences
23 (Table 1) which will cause a high degree of homoplasy. We, therefore, entertain methods for
24 estimation that explicitly account for unequal nucleotide content and multiple hits in the same
25 site such as maximum likelihood methods and the YN00 method (Yang and Nielsen 2000).
26 Yang and Nielsen (2000) (Yang and Nielsen 2000) showed that for short sequences, some

1 counting methods, such as the YN00 method, can perform better in terms of Mean Squared
2 Error (MSE) for estimating d_N and d_S . However, it is unclear in the current case how best to
3 estimate d_S . For this reason, we performed a small simulations study (see Methods) for
4 evaluating the performance of the maximum likelihood (ML) estimator of d_N and d_S (as
5 implemented in codeml (Yang 2007)) under the F3x4 model and the YN00 method implemented
6 in PAML. In general, we find that estimates under the YN00 are more biased with slightly higher
7 MSE than the ML estimate for values in the most relevant regime of $d_S < 1.5$ (Fig. 3). However,
8 we also notice that both estimators are biased under these conditions. For this reason, we
9 perform a bias correction calibrated using simulations specific to the nucleotide frequencies and
10 d_N/d_S ratio observed for SARS-CoV-2 (see Methods). The bias corrections we obtain are $\hat{d}_S^* =$
11 $\hat{d}_S + 0.455\hat{d}_S^2 - 0.824\hat{d}_S^3 + 0.264\hat{d}_S^4$, for the ML estimator and $\hat{d}_S^* = \hat{d}_S + 1.492\hat{d}_S^2 - 3.166\hat{d}_S^3 +$
12 $1.241\hat{d}_S^4$ for yn00. Notice that there is a trade-off between mean and variance (Fig. 3) so that
13 the MSE becomes very large, particularly for the for yn00 method, after bias correction. For d_S
14 > 2 the estimates are generally not reliable, however, we note that for $d_S < 1.5$ the bias-corrected
15 ML estimator tends overall to have slightly lower MSE, and we, therefore, use this estimator for
16 analyses of 300 bp regions.

17

18 *Synonymous divergence*

19 We estimate d_N and d_S under the F3x4 model in codeml (Goldman and Yang 1994; Muse and
20 Gaut 1994) and find genome-wide estimates of $d_S = 0.1604$, $d_N = 0.0065$ ($d_N/d_S = 0.0405$)
21 between SARS-CoV-2 and RaTG13 and 0.2043 ($d_N/d_S = 0.1077$) between SARS-CoV-2 and
22 RmYN02. However, a substantial amount of this divergence might be caused by recombination
23 with more divergent strains. We, therefore, infer recombination tracts using 3SEQ (Lam, et al.
24 2018) (Supplementary Table 5) and also estimate d_N and d_S for the regions with inferred
25 recombination tracts removed from all sequences (Table 3). We then find values of $d_S = 0.1462$

1 (95% C.I., 0.1340-0.1584) and $d_S = 0.1117$ (95% C.I., 0.1019-0.1215) between SARS-CoV-2
2 and RaTG13 and RmYN02, respectively. This confirms that RmYN02 is the virus most closely
3 related to SARS-CoV-2. The relative high synonymous divergence also shows that the apparent
4 high nucleotide similarity between SARS-CoV-2 and the bat strains (96.2% (Zhou, Yang, et al.
5 2020) and 97.2% (Zhou, Chen, et al. 2020)) is caused by conservation at the amino acid level
6 ($d_N/d_S = 0.0410$ and 0.0555) exacerbated by a high degree of synonymous homoplasy
7 facilitated by a highly unusual nucleotide composition.

8 The synonymous divergence to the pangolin sequences GD410721 and GX_P1E in
9 genomic regions with inferred recombination tracts removed is 0.5095 (95% C.I., 0.4794-
10 0.5396) and 1.0304 (95% C.I., 0.9669-1.0939), respectively. Values for other comparisons are
11 shown in Tables 2 and 3. In comparisons between SARS-CoV-2 and more distantly related
12 strains, d_S will be larger than 1, and with this level of saturation, estimation of divergence is
13 associated with high variance and may be highly dependent on the accuracy of the model
14 assumptions. This makes phylogenetic analyses based on synonymous mutations unreliable
15 when applied to these more divergent sequences. Nonetheless, the synonymous divergence
16 levels seem generally quite compatible with a molecular clock with a d_S of 0.9974 (95% C.I.,
17 0.9381-1.0567), 1.0366 (95% C.I., 0.9737-1.0995), 1.0333 (95% C.I., 0.9699-1.0967) and
18 1.0304 (95% C.I., 0.9669-1.0939) between the outgroup, GX_P1E, and the three ingroup
19 strains. The largest value is observed for RaTG13 (1.0366), despite this sequence being the
20 most early sampled sequence, perhaps caused by additional undetected recombination into
21 RaTG13. Another possibility is that RaTG13 has been maintained under conditions which have
22 allowed it to continue to evolve after its initial sampling.

23

24 *Sliding windows of synonymous divergence*

25 To address the issue of possible recombination we plot d_S between SARS-CoV-2, GD410721,
26 and RaTG13 and the ratio of $d_S(\text{SARS-CoV-2, GD410721})$ to $d_S(\text{SARS-CoV-2, RaTG13})$ in 300

1 bp sliding windows along the genome. Notice that we truncate the estimate of d_S at 3.0.
2 Differences between estimates larger than 2.0 should not be interpreted strongly, as these
3 estimates have high variance and likely will be quite sensitive to the specifics of the model
4 assumptions.

5 We find that $d_S(\text{SARS-CoV-2, GD410721})$ approximately equals $d_S(\text{GD410721,}$
6 $\text{RaTG13})$ and is larger than $d_S(\text{SARS-CoV-2, RaTG13})$ in almost the entire genome showing
7 than in these parts of the genome GD410721 is a proper outgroup to (SARS-CoV-2, RaTG13) .
8 One noticeable exception from this is the RBD region of the *S* gene. In this region the
9 divergence between SARS-CoV-2 and GD410721 is substantially lower than between
10 GD410721 and RaTG13 (Fig. 4A,C). The same region also has much smaller divergence
11 between SARS-CoV-2 and GD410721 than between SARS-CoV-2 and RaTG13 (Fig. 4A,C).
12 The pattern is quite different than that observed in the rest of the genome, most easily seen by
13 considering the ratio of $d_S(\text{SARS-CoV-2, GD410721})$ to $d_S(\text{SARS-CoV-2, RaTG13})$ (Fig. 2B,D).
14 In fact, the estimates of $d_S(\text{SARS-CoV-2, RaTG13})$ are saturated in this region, even though
15 they are substantially lower than 1 in the rest of the genome. This strongly suggests a
16 recombination event in the region and provides independent evidence of that previously
17 reported based on amino acid divergence (e.g., (Zhang, et al. 2020)).

18 The combined evidence from synonymous divergence and the topological recombination
19 inference procedure in 3SEQ, provide strong support for the recombination hypothesis.
20 However, these analyses alone do not distinguish between recombination into RaTG13 from an
21 unknown source as hypothesized by Boni et al. (Boni, et al. 2020) and recombination between
22 SARS-CoV-2 and GD410721 as proposed as one possible explanation by Lam et al. (Lam, et
23 al. 2020). To distinguish between these hypotheses we searched for sequences that might be
24 more closely related, in the RBD region, to RaTG13 than SARS-CoV-2 and we plotted sliding
25 window similarities across the genome for RaTG13 (Fig. 1C). We observe relatively low
26 sequence identity between RaTG13 and all three other strains in the *ACE2* contact residue

1 region of the *spike* protein, which is more consistent with the hypothesis of recombination into
2 RaTG13, as proposed by Boni et al. (Boni, et al. 2020). Moreover, our BLAST search analyses
3 of RaTG13 in this region show highest local sequence similarity with GX pangolin virus strains
4 which is the genome-wide outgroup for the three other sequences (Lam, et al. 2020) (Lam, et al.
5 2020). This observation is more compatible with the hypothesis of recombination from a virus
6 related to GX pangolin strains, than with recombination between SARS-CoV-2 and GD410721.

7 Unfortunately, because of the high level of synonymous divergence to the nearest
8 outgroup, tree estimation in small windows is extremely labile in this region. In fact, synonymous
9 divergence appears fully saturated in the comparison with GX_P1E, eliminating the possibility to
10 infer meaningful trees based on synonymous divergence. However, we can use the overall
11 maximum likelihood tree using both synonymous and nonsynonymous mutations (Fig. 2D). The
12 ML tree using sequence from the *ACE2* contact residue region supports the clustering of SARS-
13 CoV-2 and GD410721, but with unusual long external branches for all strains except SARS-
14 CoV-2, possibly reflecting smaller recombination regions within the *ACE2* contact residue
15 region.

16

17 *Weakly deleterious mutations and clock calibrations*

18 The use of synonymous mutations provides an opportunity to calibrate the molecular clock
19 without relying on amino acid changing mutations that are more likely to be affected by
20 selection. The rate of substitution of weakly and slightly deleterious mutations is highly
21 dependent on ecological factors and the effective population size. Weakly deleterious mutations
22 are more likely to be observed over small time scales than over long time scales, as they are
23 unlikely to persist in the population for a long time and go to fixation. This will lead to a
24 decreasing d_N/d_S ratio for longer evolutionary lineages. Furthermore, changes in effective
25 population size will translate into changes in the rate of substitution of slightly deleterious
26 mutations. Finally, changes in ecology (such as host shifts, host immune changes, changes in

1 cell surface receptor, etc.) can lead to changes in the rate of amino acid substitution. For all of
2 these reasons, the use of synonymous mutations, which are less likely to be the subject of
3 selection than nonsynonymous mutations, are preferred in molecular clock calculations. For
4 many viruses, the use of synonymous mutations to calibrate divergence times is not possible,
5 as synonymous sites are fully saturated even at short divergence times. However, for the
6 comparisons between SARS-CoV-2 and RaTG13, and SARS-CoV-2 and RmYN02,
7 synonymous sites are not saturated and can be used for calibration. We find an estimate of $\omega =$
8 0.0391 between SARS-CoV-2 and RaTG13, excluding just the small RDB region inferred by
9 3SEQ to be recombined in SARS-CoV-2 (Supplementary Table 5, coordinates: 22851-23094).
10 Using 1000 parametric simulations under the estimated values and the F3x4 codon model, we
11 find that the estimate is approximately unbiased ($\hat{\omega} = 0.0398$, S.E.M.= 0.0001) and with
12 standard deviation 0.0033, providing an approximate 95% confidence interval of (0.0332,
13 0.0464). Also, using 59 human strains of SARS-CoV-2 from Genbank and National Microbiology
14 Data Center (See Methods) we obtain an estimate of $\omega = 0.5604$ using the F3x4 model in
15 codeml. Notice that there is a 14-fold difference in d_N/d_S ratio between these estimates.
16 Assuming very little of this difference is caused by positive selection, this suggests that the vast
17 majority of mutations currently segregating in the SARS-CoV-2 are slightly or weakly deleterious
18 for the virus.

19

20 *Dating of divergence between Bat viruses and SARS-CoV-2*

21 To calibrate the clock we use the estimate provided by ([http://virological.org/t/phylogenetic-
22 analysis-of-sars-cov-2-update-2020-03-06/420](http://virological.org/t/phylogenetic-analysis-of-sars-cov-2-update-2020-03-06/420)) of $\mu = 1.04 \times 10^{-3}$ substitutions/site/year (95% CI:
23 0.71×10^{-3} , 1.40×10^{-3}). The synonymous specific mutation rate can be found from this as
24 $d_S/\text{year} = \mu_S = \mu/(pS + \omega pN)$, where pN and pS are the proportions of nonsynonymous and
25 synonymous sites, respectively. The estimate of the total divergence on the two lineages is then
26 $\hat{t} = dS(pS + \omega pN)/\mu$. Inserting the numbers from Table 3 for the divergence between SARS-

1 CoV-2 and RaTG13 and RmYN02 ,respectively, we find a total divergence of 96.92 years and
2 74.05 years respectively. Taking into account that RaTG13 was isolated July 2013, we find an
3 estimated tMRCA between that strain and SARS-CoV-2 of $\hat{t} = (96.92 + 6.5)/2 = 51.71$ years.
4 Similarly, we find an estimate of divergence between SARS-CoV-2 and RmYN02 of $\hat{t} = 74.05/2$
5 = 37.02 years, assuming approximately equal sampling times. The estimate for SARS-CoV-2
6 and RaTG13 is compatible with the values obtained by Boni *et al.* (Boni, et al. 2020), who used
7 quite different methods for dating. The variance in the estimate in d_S is small and the uncertainty
8 is mostly dominated by the uncertainty in the estimate of the mutation rate. We estimate the
9 S.D. in \hat{t} using 1000 parametric simulations, using the ML estimates of all parameters, for both
10 RaTG13 vs. SARS-CoV-2 and for RmYN02 vs. SARS-CoV-2, and for each simulated data also
11 simulating values of μ and ω from normal distributions with mean 1.04×10^{-3} and S.D. 0.18×10^{-3} ,
12 and mean 0.5604 and S.D. 0.1122, respectively. We subject each simulated data set to the
13 same inference procedure as done on the real data. Our estimate of the S.D. in the estimate is
14 11.8 for RaTG13 vs. SARS-CoV-2 and 9.41 for RmYN02 vs. SARS-CoV-2, providing an
15 approximate 95% confidence interval of (28.11, 75.31) and (18.19, 55.85), respectively. For
16 RaTG13, if including all sites, except the 244-bp in the RBD of the S gene (Supplementary
17 Table 5), the estimate is 55.02 years with an approx. 95% C.I. of (29.4, 80.7). As more SARS-
18 CoV-2 sequences are being obtained, providing more precise estimates of the mutation rate,
19 this confidence interval will become narrower. However, we warn that the estimate is based on
20 a molecular clock assumption and that violations of this assumption eventually will become a
21 more likely source of error than the statistical uncertainty quantified in the calculation of the
22 confidence intervals. We also note that, so far, we have assumed no variation in the mutation
23 rate among synonymous sites. However, just from the analysis of the 300 bp windows, it is clear
24 that is not true. The variance in the estimate of d_S among 300 bp windows from the RaTG13-
25 SARS-CoV-2 comparison is approximately 0.0113. In contrast, in the simulated data assuming
26 constant mutation rate, the variance is approximately 0.0034, suggesting substantial variation in

1 the synonymous mutation rate along the length of the genome. Alternatively, this might be
2 explained by undetected recombination in the evolutionary history since the divergence of the
3 strains.

4

5 **Discussion**

6 The unusually skewed distribution of nucleotide frequencies in synonymous sites in SARS-CoV-
7 2 (Kandeel, et al. 2020), along with high divergence, complicates the estimation of synonymous
8 divergence in SARS-CoV-2 and related viruses. In particular, in the third codon position the
9 nucleotide frequency of T is 43.5% while it is just 15.7% for C. This resulting codon usage is not
10 optimized for mammalian cells (e.g, Chamary, et al. 2006). A possible explanation is a strong
11 mutational bias caused by Apolipoprotein B mRNA-editing enzymes (APOBECs) which can
12 cause Cytosine-to-Uracil changes (Giorgio, et al. 2020).

13 A consequence of the skewed nucleotide frequencies is a high degree of homoplasy in
14 synonymous sites that challenges estimates of d_S . We here evaluated estimators of d_S in 300 bp
15 sliding windows and found that a bias-corrected version of the maximum likelihood estimator
16 tended to perform best for values of $d_S < 2$. We used this estimator to investigate the
17 relationship between SARS-CoV-2 and related viruses in sliding windows. We show that
18 synonymous mutations show shorter divergence to pangolin viruses, than the otherwise most
19 closely related bat virus, RaTG13, in part of the receptor-binding domain of the *spike* protein.
20 This strongly suggests that the previously reported amino acid similarity between pangolin
21 viruses and SARS-CoV-2 is not due to convergent evolution, but more likely is due to
22 recombination. Boni, et al. 2020 used the program 3SEQ to infer recombination events in
23 SARS-CoV-2 and related viruses and did not recover strong evidence for recombination specific
24 to SARS-CoV-2. In contrast in our analyses, when applying the program to a smaller subset of
25 viral strains, 3SEQ identifies recombination from pangolin strains into SARS-CoV-2. Possibly,
26 the inclusion of more divergent strains in the more comprehensive analysis by (Boni, et al.

1 2020) has lead to reduced power to detect recombination among the less divergent strains
2 closely related to SARS-CoV-2. In any case, the 3SEQ analyses presented provides further
3 support for the recombination hypothesis. However, we also find that the synonymous
4 divergence between SARS-CoV-2 and pangolin viruses in this region is relatively high,
5 suggesting that the recombination was into RaTG13 from an unknown strain, rather than
6 between pangolin viruses and SARS-CoV-2, which is consistent with the hypothesis proposed
7 by Boni, et al. 2020, and we can exclude the possibility of a very recent recombination event
8 between pangolins and SARS-CoV-2. Another alternative explanation is that there is a
9 mutational hotspot located in the RBD region. While it is possible that selection may have
10 favored hypermutability in this domain, the most parsimonious explanation at the moment, in the
11 absence of knowledge of molecular mechanisms to explain such a hotspot, is recombination
12 into the RaTG13 strain. To fully distinguish between these hypotheses, additional strains would
13 have to be discovered that either are candidates for introgression into RaTG13 or can break up
14 the lineage in the phylogenetic tree between pangolin viruses and RaTG13.

15 The fact that synonymous divergence to the outgroups, RaTG13 and RmYN02, is not
16 fully saturated, provides an opportunity for a number of different analyses. First, we can date the
17 time of the divergence between the bat viruses and SARS-CoV-2 using synonymous mutations
18 alone. In doing so, we find estimates of 51.71 years (95% C.I., 28.11-75.31) and 37.02 years
19 (95% C.I., 18.19-55.85), respectively. Most of the uncertainty in these estimates comes from
20 uncertainty in the estimate of the mutation rate reported for SARS-CoV-2. As more data is being
21 produced for SARS-CoV-2, the estimate should become more precise and the confidence
22 interval significantly narrowed. However, we warn that a residual cause of unmodeled statistical
23 uncertainty is deviations from the molecular clock. Variation in the molecular clock could be
24 modeled statistically (see e.g., Drummond, et al. 2006 and Lartillot, et al. 2016), but the fact that
25 synonymous mutations are mostly saturated for more divergent viruses that would be needed to
26 train such models, is a challenge to such efforts. On the positive side, we note that the

1 estimates of d_S given in Table 3 in general are highly compatible with a constant molecular
2 clock.

3 Another advantage of estimation of synonymous and nonsynonymous rates in the
4 outgroup lineage, is that it can provide estimates of the mutational load of the current pandemic.
5 The d_N/d_S ratio is almost 14 times larger in the circulating SARS-CoV-2 strains than in the
6 outgroup lineage. While some of this difference could possibly be explained by positive
7 selection acting at a higher rate after zoonotic transfer, it is perhaps more likely that a
8 substantial proportion of segregating nonsynonymous mutations are deleterious, suggesting a
9 very high and increasing mutation load in circulating SARS-CoV-2 strains.

10

11 **Acknowledgements**

12 We are grateful to Dr. Yongyi Shen and Dr. E.C Holmes for providing the genome sequence of
13 GD410721 and RmYN02, respectively. We also thank Dr. Adi Stern for discussion. The
14 research was funded by Koret-UC Berkeley-Tel Aviv University Initiative in Computational
15 Biology and Bioinformatics to RN.

16

17 **Reference**

18 Andersen KG, Rambaut A, Lipkin WI, Holmes EC, Garry RF. 2020. The proximal origin of SARS-CoV-2.
19 Nat Med 26:450-452.

20 Boni MF, Lemey P, Jiang X, Lam TT, Perry B, Castoe T, Rambaut A, Robertson DL. 2020. Evolutionary
21 origins of the SARS-CoV-2 sarbecovirus lineage responsible for the COVID-19 pandemic. bioRxiv
22 <https://doi.org/10.1101/2020.03.30.015008>.

23 Chamary JV, Parmley JL, Hurst LD. 2006. Hearing silence: non-neutral evolution at synonymous sites in
24 mammals. Nat Rev Genet 7:98-108.

25 Drummond AJ, Ho SY, Phillips MJ, Rambaut A. 2006. Relaxed phylogenetics and dating with confidence.
26 PLoS Biol 4:e88.

- 1 Felsenstein J. 2009. PHYLIP (Phylogeny Inference Package) version 3.7a. Distributed by the author.
2 Department of Genome Sciences, University of Washington, Seattle.
- 3 Giorgio SD, Martignano F, Torcia MG, Mattiuz G, Conticello SG. 2020. Evidence for RNA editing in the
4 transcriptome of 2019 Novel Coronavirus. bioRxiv <https://doi.org/10.1101/2020.03.02.973255>.
- 5 Goldman N, Yang Z. 1994. A codon-based model of nucleotide substitution for protein-coding DNA
6 sequences. *Mol Biol Evol* 11:725-736.
- 7 Hu D, Zhu C, Ai L, He T, Wang Y, Ye F, Yang L, Ding C, Zhu X, Lv R, et al. 2018. Genomic
8 characterization and infectivity of a novel SARS-like coronavirus in Chinese bats. *Emerg Microbes Infect*
9 7:154.
- 10 Kandeel M, Ibrahim A, Fayez M, Al-Nazawi M. 2020. From SARS and MERS CoVs to SARS-CoV-2:
11 Moving toward more biased codon usage in viral structural and nonstructural genes. *J Med Virol*.
- 12 Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7: improvements in
13 performance and usability. *Mol Biol Evol* 30:772-780.
- 14 Lam HM, Ratmann O, Boni MF. 2018. Improved Algorithmic Complexity for the 3SEQ Recombination
15 Detection Algorithm. *Mol Biol Evol* 35:247-251.
- 16 Lam TT, Shum MH, Zhu HC, Tong YG, Ni XB, Liao YS, Wei W, Cheung WY, Li WJ, Li LF, et al. 2020.
17 Identifying SARS-CoV-2 related coronaviruses in Malayan pangolins. *Nature*.
- 18 Lartillot N, Phillips MJ, Ronquist F. 2016. A mixed relaxed clock model. *Philos Trans R Soc Lond B Biol*
19 *Sci* 371.
- 20 Li Q, Guan X, Wu P, Wang X, Zhou L, Tong Y, Ren R, Leung KSM, Lau EHY, Wong JY, et al. 2020. Early
21 Transmission Dynamics in Wuhan, China, of Novel Coronavirus-Infected Pneumonia. *N Engl J Med*
22 382:1199-1207.
- 23 Loytynoja A. 2014. Phylogeny-aware alignment with PRANK. *Methods Mol Biol* 1079:155-170.
- 24 Lu R, Zhao X, Li J, Niu P, Yang B, Wu H, Wang W, Song H, Huang B, Zhu N, et al. 2020. Genomic
25 characterisation and epidemiology of 2019 novel coronavirus: implications for virus origins and receptor
26 binding. *Lancet* 395:565-574.
- 27 Martin DP, Murrell B, Golden M, Khoosal A, Muhire B. 2015. RDP4: Detection and analysis of
28 recombination patterns in virus genomes. *Virus Evol* 1:vev003.
- 29 Muse SV, Gaut BS. 1994. A likelihood approach for comparing synonymous and nonsynonymous
30 nucleotide substitution rates, with application to the chloroplast genome. *Mol Biol Evol* 11:715-724.

- 1 Nguyen LT, Schmidt HA, von Haeseler A, Minh BQ. 2015. IQ-TREE: a fast and effective stochastic
2 algorithm for estimating maximum-likelihood phylogenies. *Mol Biol Evol* 32:268-274.
- 3 Patiño-Galindo JÁ, Filip I, AlQuraishi M, Rabadan R. 2020. Recombination and lineage-specific mutations
4 led to the emergence of SARS-CoV-2. bioRxiv DOI: 10.1101/2020.02.10.942748.
- 5 Wan Y, Shang J, Graham R, Baric RS, Li F. 2020. Receptor Recognition by the Novel Coronavirus from
6 Wuhan: an Analysis Based on Decade-Long Structural Studies of SARS Coronavirus. *J Virol* 94.
- 7 Wu F, Zhao S, Yu B, Chen YM, Wang W, Song ZG, Hu Y, Tao ZW, Tian JH, Pei YY, et al. 2020. A new
8 coronavirus associated with human respiratory disease in China. *Nature* 579:265-269.
- 9 Xiao K, Zhai J, Feng Y, Zhou N, Zhang X, Zou J-J, Li N, Guo Y, Li X, Shen X, et al. 2020. Isolation and
10 Characterization of 2019-nCoV-like Coronavirus from Malayan Pangolins bioRxiv
11 <https://doi.org/10.1101/2020.02.17.951335>.
- 12 Yang Z. 2007. PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol* 24:1586-1591.
- 13 Yang Z, Nielsen R. 2000. Estimating synonymous and nonsynonymous substitution rates under realistic
14 evolutionary models. *Mol Biol Evol* 17:32-43.
- 15 Yang Z, Nielsen R, Goldman N, Pedersen AM. 2000. Codon-substitution models for heterogeneous
16 selection pressure at amino acid sites. *Genetics* 155:431-449.
- 17 Zhang T, Wu Q, Zhang Z. 2020. Probable Pangolin Origin of SARS-CoV-2 Associated with the COVID-19
18 Outbreak. *Curr Biol* 30:1346-1351 e1342.
- 19 Zhang YZ, Holmes EC. 2020. A Genomic Perspective on the Origin and Emergence of SARS-CoV-2. *Cell*
20 181:223-227.
- 21 Zhou H, Chen X, Hu T, Li J, Song H, Liu Y, Wang P, Liu D, Yang J, Holmes EC, et al. 2020. A novel bat
22 coronavirus reveals natural insertions at the S1/S2 cleavage site of the Spike protein and a possible
23 recombinant origin of HCoV-19. bioRxiv doi: <https://doi.org/10.1101/2020.03.02.974139>.
- 24 Zhou P, Yang XL, Wang XG, Hu B, Zhang L, Zhang W, Si HR, Zhu Y, Li B, Huang CL, et al. 2020. A
25 pneumonia outbreak associated with a new coronavirus of probable bat origin. *Nature* 579:270-273.
- 26
27
28

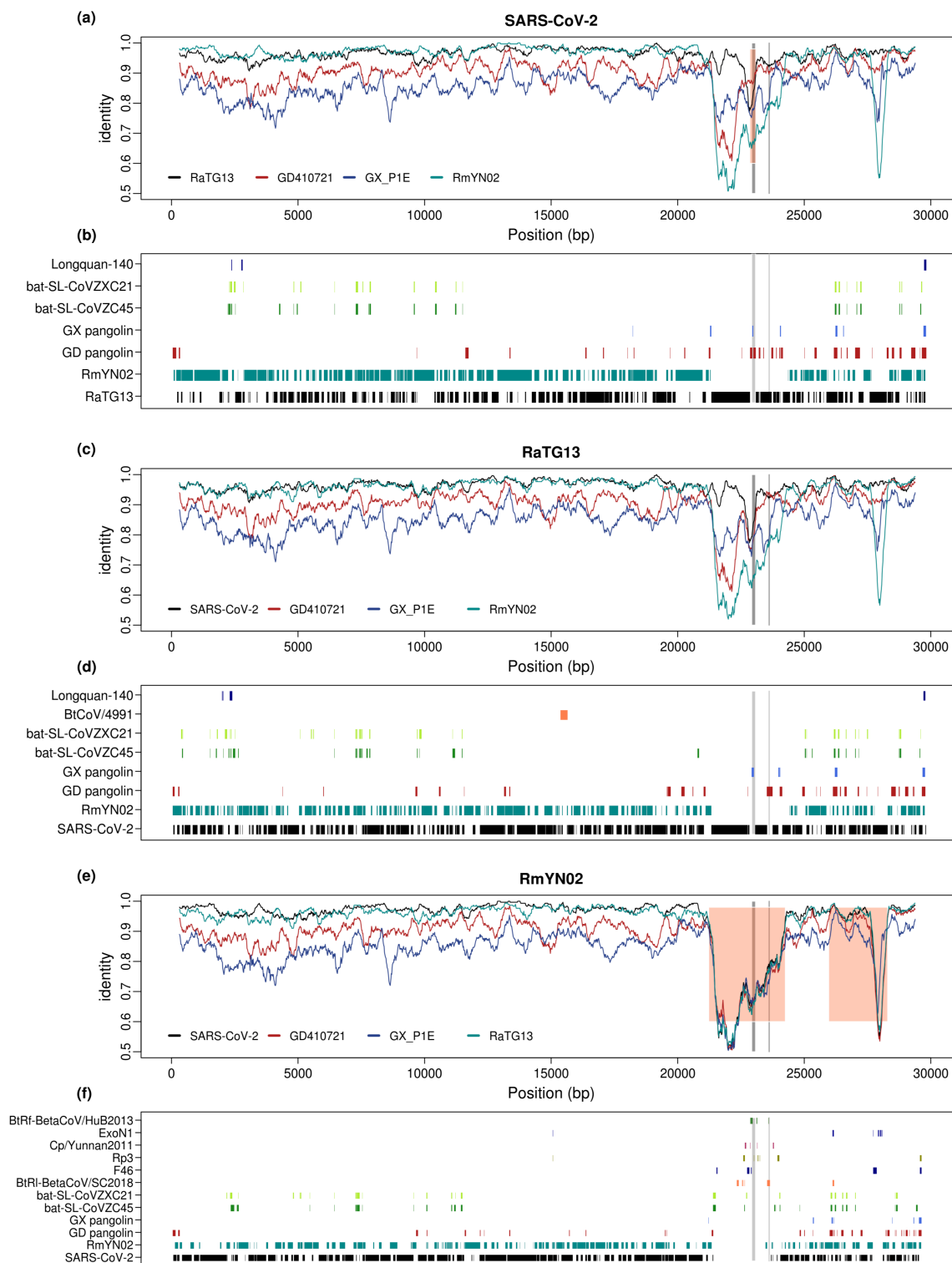


Figure 1. Genome-wide identity plot and top blast hits for SARS-CoV-2, RaTG13 and RmYN02. (a) 300 bp sliding-windows of nucleotide identity between SARS-CoV-2 and the four most closely related viral strains, RmYN02, RaTG13, GD410721 and GX_P1E. Orange shading

marks the recombinant region in SARS-CoV-2 inferred by 3SEQ (details in Supplementary Table 5). (b) the plot lists all the viral strains that are the unique best BLAST hit in at least three 100-bp windows, when blasting with SARS-CoV-2, with the regions where each strain is the top blast hit marked. (b) and (c). Figures for RaTG13 (c, d) and RmYN02 (e, f) generated in the same way as for SARS-CoV-2 in (a) and (b). The ACE2 contact residues of RBD region (left) and the furin sites (right) of the S protein are marked in both plots with grey lines.

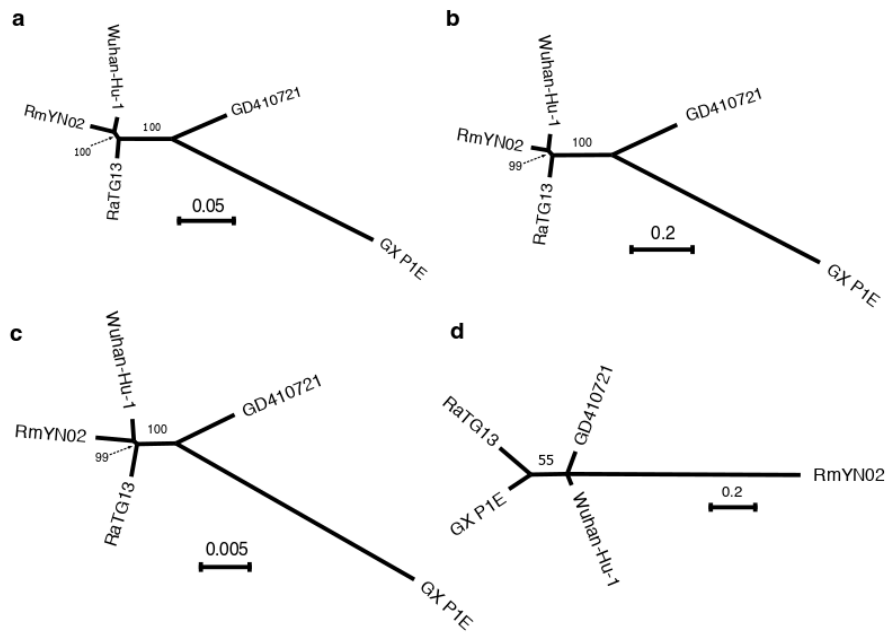


Figure 2. Unrooted phylogenies of the virus strains. (a) Maximum-likelihood tree in genomic regions with recombination tracts removed; (b) Neighbor-joining tree using synonymous mutation (d_S) distance in genomic regions with recombination tracts removed; (c) Neighbor-joining tree using non-synonymous mutation (d_N) distances in genomic regions with recombination tracts removed; (d) The maximum-likelihoods tree at the receptor-binding domain ACE2 contact residues (51 amino acids) region. The bootstrap values are based on 1,000 replicates. The associated distance matrix for (b) and (c) can be found in Table 3.

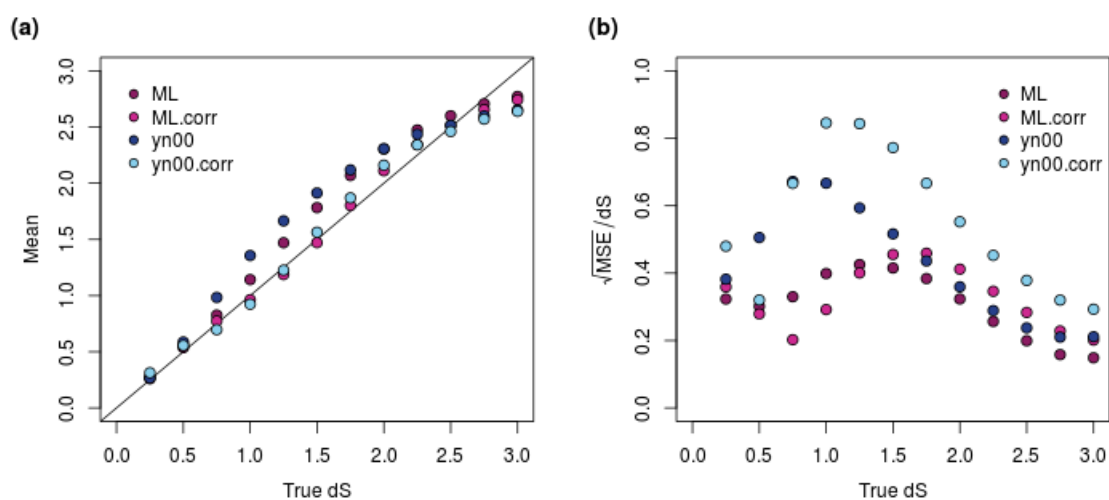


Figure 3. Bias correction for d_S estimate in 300-bp windows. (a) The mean of d_S estimates using different methods; ML.corr and yn00.corr are the bias corrected versions of the ML and yn00 methods, respectively. (b) Errors in d_S estimates as measured using the ratio of square root of mean squared error (MSE) to true d_S . All the estimates are based on 10,000 simulations. ML: maximum-likelihood estimates using the f3x4 model in codeml; ML.corr, maximum-likelihood estimates with bias correction; yn00, count-based estimates in (Yang and Nielsen 2000); yn00.corr, yn00 estimates with bias correction. All d_S estimates are truncated at 3, explaining the reduction in MSE with increasing values of d_S as d_S approaches 3.

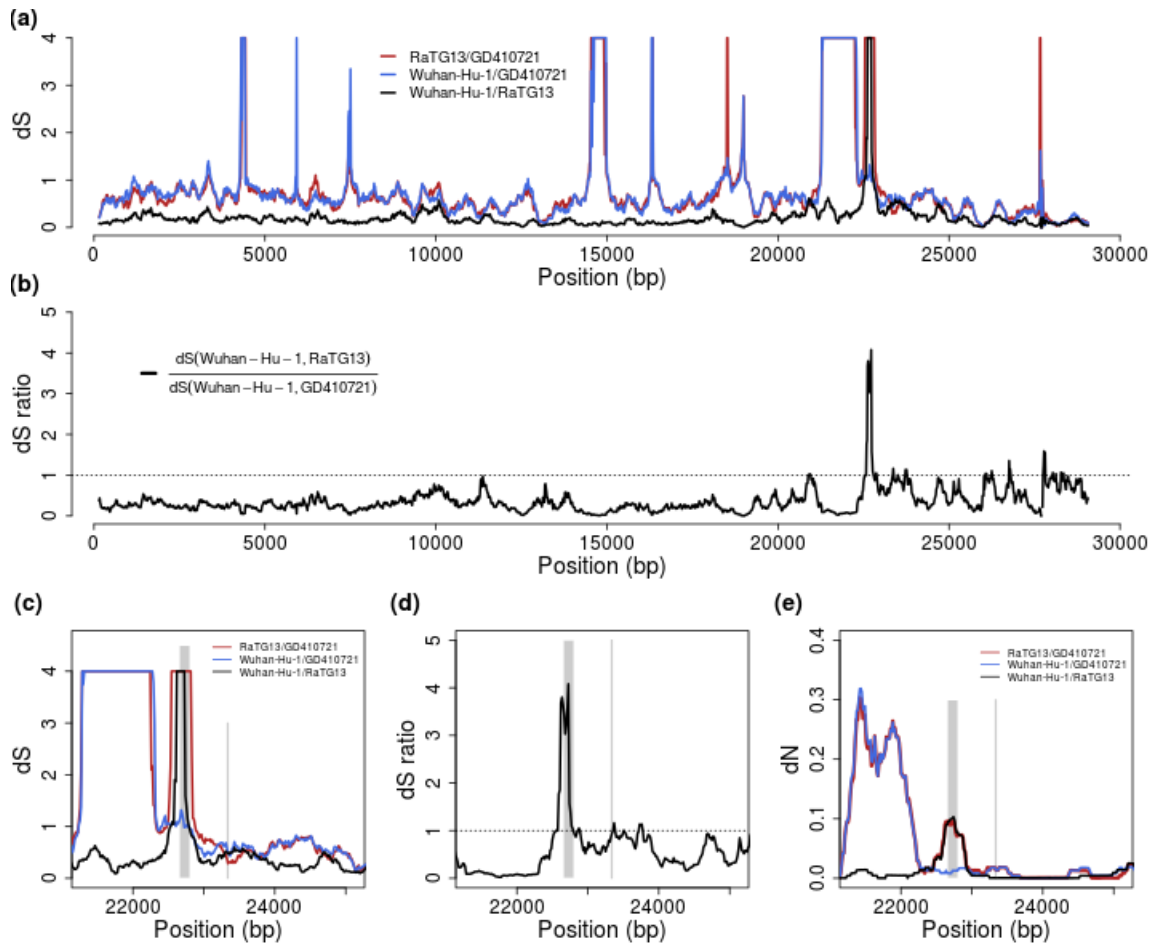


Figure 4. d_S and d_N estimates across the virus genome. (a) Pairwise d_S estimates in 300-bp sliding windows for RaTG13, GD410721 and Wuhan-Hu-1, the estimates are truncated at 4. (b) d_S ratio of $d_S(\text{Wuhan-Hu-1, RaTG13})$ to $d_S(\text{Wuhan-Hu-1, GD410721})$. (c) and (d) are the zoom-in plot for d_S and d_S -ratio at the *spike* (S) protein region. The receptor-binding domain contact residues (left) and furin site regions (right) are marked with grey lines. (e) the pairwise d_N estimates in 300-bp sliding windows in the S protein for these strains. The d_S values are truncated at 4 in the plots.

Table 1. Genome-wide nucleotide composition at the third position of the codons in the viral strains. The nucleotide compositions at the first and second positions can be found in Supplementary table 1.

Accession	T	C	A	G
GD410721	42.71%	16.17%	28.55%	12.57%
GX_P1E	42.52%	16.40%	28.27%	12.81%
RaTG13	43.57%	15.74%	27.98%	12.71%
RmYN02	43.31%	15.90%	27.98%	12.81%
Wuhan-Hu-1	43.49%	15.73%	28.16%	12.62%

Table 2. Whole genome d_N and d_S estimates among the viral strains. The d_S estimates are shaded in green, and the d_N estimates are in orange shade. The 95% confidence intervals, calculated based on 1,000 bootstrap replicates, are included in the brackets for each estimates.

	GD410721	GX_P1E	RaTG13	RmYN02	Wuhan-Hu-1
GD410721		0.0372 (0.0341-0.0403)	0.0171 (0.0152-0.0190)	0.0293 (0.0266-0.0320)	0.0160 (0.0142-0.0178)
GX_P1E	0.9883 (0.9338-1.0428)		0.0347 (0.0318-0.0376)	0.0485 (0.0450-0.0520)	0.0342 (0.0314-0.0370)
RaTG13	0.5392 (0.5105-0.5679)	1.0156 (0.9608-1.0704)		0.0235 (0.0210-0.0260)	0.0065 (0.0053-0.0077)
RmYN02	0.6001 (0.5681-0.6321)	1.0757 (1.0166-1.1348)	0.2438 (0.2285-0.2591)		0.0220 (0.0195-0.0245)
Wuhan-Hu-1	0.5425 (0.5131-0.5719)	0.9973 (0.9434-1.0512)	0.1604 (0.1491-0.1717)	0.2043 (0.1901-0.2185)	

Table 3. Genome-wide d_N and d_S estimates after removing recombination regions inferred by 3SEQ. The d_S estimates are shaded in green, and the d_N estimates are in orange shade. The coordinates relative to the Wuhan-Hu-1 genome of the masked region can be found in the method section. The 95% confidence intervals, calculated based on 1,000 bootstrap replicates, are included in the brackets for each estimates.

	GD410721	GX_P1E	RaTG13	RmYN02	Wuhan-Hu-1
GD410721		0.0348 (0.0317-0.0379)	0.0138 (0.0120-0.0156)	0.0152 (0.0133-0.0171)	0.0135 (0.0117-0.0153)
GX_P1E	0.9974 (0.9381-1.0567)		0.0357 (0.0325-0.0389)	0.0361 (0.0329-0.0393)	0.0349 (0.0318-0.0380)
RaTG13	0.4962 (0.4669-0.5255)	1.0366 (0.9737-1.0995)		0.0079 (0.0066-0.0092)	0.0060 (0.0048-0.0071)
RmYN02	0.5070 (0.4773-0.5366)	1.0333 (0.9699-1.0967)	0.1522 (0.1395-0.1649)		0.0062 (0.0050-0.0074)
Wuhan-Hu-1	0.5095 (0.4794-0.5396)	1.0304 (0.9669-1.0939)	0.1462 (0.1340-0.1584)	0.1117 (0.1019-0.1215)	