

1 **TITLE PAGE**

2 **Title:** THE ORIGIN OF A NEW HUMAN VIRUS: PHYLOGENETIC ANALYSIS OF THE
3 EVOLUTION OF SARS-COV-2

4

5 **RUNNING TITLE:** Phylogenetic analysis and evolution of SARS-CoV-2

6

7 **Authors:**

8 Matías J. PERESON^{a,b}, Laura MOJSIEJCZUK^{a,b}, Alfredo P. MARTÍNEZ^c, Diego M.
9 FLICHMAN^{b,c}, Gabriel H. GARCIA^a, Federico A. DI LELLO^{a,b}

10

11 **Affiliations:**

12 ^aUniversidad de Buenos Aires. Facultad de Farmacia y Bioquímica. Instituto de
13 Investigaciones en Bacteriología y Virología Molecular (IBaViM). Buenos Aires, Argentina.

14 ^bConsejo Nacional de Investigaciones Científicas y Técnicas (CONICET), Ciudad Autónoma
15 de Buenos Aires, Argentina.

16 ^cVirology Section, Centro de Educación Médica e Investigaciones Clínicas Norberto Quirno
17 "CEMIC". Buenos Aires, Argentina.

18 ^dInstituto de Investigaciones Biomédicas en Retrovirus y Síndrome de Inmunodeficiencia
19 Adquirida (INBIRS) – Consejo Nacional de Investigaciones Científicas y Técnicas
20 (CONICET), Universidad de Buenos Aires, Buenos Aires, Argentina.

21

22 **Disclosures:** None

23

24 **Corresponding author:**

25 Dr. Federico Alejandro Di Lello, Facultad de Farmacia y Bioquímica, Universidad de Buenos
26 Aires, Instituto de Investigaciones en Bacteriología y Virología Molecular (IBaViM). Junín 956,
27 4º piso, (1113), Ciudad Autónoma de Buenos Aires, Argentina.
28 Phone: +54 11 5287 4472, Fax: +54 11 5287 4662, E-mail: fadilello@ffyb.uba.ar

29

30 **Author contributions**

31

32 MJP: Data curation, acquisition of data, analysis and interpretation of data, drafting the article,
33 final approval of the version to be submitted.

34

35 LM: Data curation, acquisition of data, analysis and interpretation of data, revising the article
36 critically for important intellectual content, final approval of the version to be submitted.

37

38 APM: Data curation, Validation, revising the article critically for important intellectual content,
39 final approval of the version to be submitted.

40

41 DMF: Data curation, Validation, drafting the article, final approval of the version to be
42 submitted.

43

44 GG: Data curation, acquisition of data, analysis and interpretation of data, drafting the article,
45 final approval of the version to be submitted.

46

47 FAD: Conception and design of the study, acquisition of data, analysis and interpretation of
48 data, drafting the article, final approval of the version to be submitted.

49

50

51

52

53

54

55 **ABSTRACT**

56 *Objectives:* During the first months of SARS-CoV-2 evolution in a new host, contrasting
57 hypotheses have been proposed about the way the virus has evolved and diversified
58 worldwide. The aim of this study was to perform a comprehensive evolutionary analysis to
59 describe the human outbreak and the evolutionary rate of different genomic regions of SARS-
60 CoV-2.

61 *Methods:* The molecular evolution in nine genomic regions of SARS-CoV-2 was analyzed
62 using three different approaches: phylogenetic signal assessment, emergence of amino acid
63 substitutions, and Bayesian evolutionary rate estimation in eight successive fortnights since
64 the virus emergence.

65 *Results:* All observed phylogenetic signals were very low and consistent trees were obtained.
66 However, after four months of evolution, it was possible to identify regions revealing an
67 incipient viral lineages formation despite the low phylogenetic signal, since fortnight 3. Finally,
68 the SARS-CoV-2 evolutionary rate for regions nsp3 and S, the ones presenting greater
69 variability, was estimated to range from 1.37 to 2.19×10^{-3} substitution/site/year.

70 *Conclusions:* In conclusion, results obtained in this work about the variable diversity of crucial
71 viral regions and the determination of the evolutionary rate are consequently decisive to
72 understanding essential feature of viral emergence. In turn, findings may allow identifying the
73 best targets for antiviral treatments and vaccines development.

74

75 **KEYWORDS:** SARS-CoV-2, Phylogeny, Evolution, Evolutionary Rate

76 **1. Introduction**

77 Coronaviruses belong to Coronaviridae family and have a single strand of positive-sense RNA
78 genome of 26 to 32 kb in length (Su et al. 2016). They have been identified in different avian
79 hosts (Cavanagh, 2007, Ismail et al. 2003) as well as in various mammals including bats,
80 mice, dogs, etc. Periodically, new mammalian coronaviruses are identified. In late December
81 2019, Chinese health authorities identified groups of patients with pneumonia of unknown
82 cause in Wuhan, Hubei Province, China (Zhu et al. 2020). The pathogen, a new coronavirus
83 called SARS-CoV-2 (Coronaviridae Study Group of the International Committee on Taxonomy
84 of Viruses, 2020), was identified by local hospitals using a surveillance mechanism for
85 "pneumonia of unknown etiology" (Li et al. 2020a, Li et al. 2020b, Zhu et al. 2020). The
86 pandemic has spread rapidly and, to date, more than 14 million confirmed cases and nearly
87 600,000 deaths have been reported in just over a six months period (World Health
88 Organization, 2020). This rapid viral spread raises interesting questions about the way its
89 evolution is driven during the pandemic. From the SARS-CoV-2 genome, 16 non-structural
90 proteins (nsp1-16), 4 structural proteins [spike (S), envelope (E), membrane (M) and
91 nucleoprotein (N)], and other proteins essential to complete the replication cycle are translated
92 (Cui et al. 2019, Luk et al. 2019). The large amount of information currently available allows
93 knowing, as never before, the real-time evolution history of a virus since its interspecies jump
94 (Zhou et al. 2020). Most studies published to date have characterized the viral genome and
95 evolution by analyzing a small number of sequences (Benvenuto et al. 2020, Cagliani et al.
96 2020, Phan, 2020) since processing of complete genomes constitutes an enormous demand
97 of time and resources. Despite this, until now, the viral genomic region providing the most
98 accurate information to characterize SARS-CoV-2, could not be established. This lack of
99 information prevent from investigating its molecular evolution and monitoring biological
100 features affecting the development of antiviral and vaccines. Therefore, the aim of this study

101 was to perform a comprehensive viral evolutionary analysis in order to describe the human
102 outbreak and the molecular evolution rate of different genomic regions of SARS-CoV-2.

103 **2. Materials and Methods**

104 *2.1 Datasets*

105 In order to generate a dataset representing different geographic regions and time evolution of
106 the SARS-CoV-2 pandemic from December 24th 2019 to April 17th 2020, data of all the
107 complete genome sequences available at GISAID (<https://www.gisaid.org/>) on April 18, 2020
108 were collected. Data inclusion criteria were: a.- complete genomes, b.- high coverage level,
109 and c.- human hosts only (no other animals, cell culture, or environmental samples). Complete
110 genomes were aligned using MAFFT against the Wuhan-Hu-1 reference genome
111 (NC_045512.2, EPI_ISL_402125). The resulting multiple sequence alignment (dataset 1) was
112 split in nine datasets corresponding to nine coding regions: a.- four structural proteins
113 [envelope (E), nucleocapsid (N), spike (S), Orf3a], b.- four nonstructural proteins (nsp1, nsp3,
114 Orf6, and nsp14), and c.- an unknown function protein (Orf8).

115 More than six thousand SARS-CoV-2 publicly available nucleotide sequences were
116 downloaded. After data selection according to the inclusion criteria, 1616 SARS-CoV-2
117 complete genomes were included in dataset 1. Sequences of this dataset 1 came from 55
118 countries belonging to the five continents as follow: Africa: 39 sequences, Americas: 383
119 sequences, Asia: 387 sequences, Europe: 686 sequences and Oceania: 121 sequences. After
120 elimination of sequences with indeterminate or ambiguous positions, the number of analyzed
121 sequences for each region were: nsp1, 1608; nsp3, 1511; nsp14, 1550; S, 1488; Orf3a, 1600;
122 E, 1615; Orf6, 1616; Orf8, 1612; and N, 1610. Finally, nucleotide sequences were grouped by
123 fortnight (FN) according to their collection date. Table 1 summarizes the number of sequences
124 per fortnight since the beginning of the pandemic up to FN 8. Dataset 2 was created with the
125 variable sequences of each region analyzed in Dataset 1. Dataset 2 was used for the
126 phylogenetic signal analysis and the Bayesian coalescent trees construction.

127

128 *2.2 Phylogenetic signal*

129 To determine the phylogenetic signal of each of the nine generated alignments, Likelihood
130 Mapping analyzes were carried out (Strimmer & von Haeseler, 1997), using the Tree Puzzle
131 v5.3 program (Schmidt et al. 2002) and the Quartet puzzling algorithm. This algorithm allowed
132 analyzing the tree topologies that can be completely solved from all possible quartets of the n
133 alignment sequences using maximum likelihood. An alignment with defined tree values
134 greater than 70-80% presents strong support from the statistical point of view (Schmidt et al.
135 2002). Identical sequences were also removed with ElimDupes (Available at
136 <https://www.hiv.lanl.gov/content/sequence/elimdupesv2/elimdupes.html>) as they increase
137 computation time and provide no additional information about data phylogeny. The best-fit
138 evolutionary model to each dataset was selected based on the Bayesian Information Criterion
139 obtained with the JModelTest v2.1.10 software (Darriba et al. 2012).

140

141 *2.3 Analysis of amino acid substitutions*

142 Entropy-One (Available at
143 https://www.hiv.lanl.gov/content/sequence/ENTROPY/entropy_one.html) was used to
144 determining the frequency of amino acids at each position for the nine genomic regions
145 analyzed and evaluating their permanence in the eight investigated fortnights.

146

147 *2.4 Bayesian coalescence and phylogenetic analysis*

148 To study the relationship between SARS-CoV-2 sequences, nine regions of the virus genome
149 were investigated by Bayesian analyses. Phylogenetic trees were constructed using Bayesian
150 inference with MrBayes v3.2.7a (Ronquist et al. 2012). Each gene was analyzed
151 independently with the same dataset used for the phylogenetic signal analysis so that non-
152 identical sequences were included in the analysis. Analyses were run for five million

153 generations and sampled every 5000 generations. Convergence of parameters [effective
154 sample size (ESS) ≥ 200 , with a 10% burn-in] was verified with Tracer v1.7.1 (Rambaut et al.
155 2018). Phylogenetic trees were visualized with FigTree v1.4.4.

156

157 *2.5 Evolutionary rate*

158 The estimation of the nucleotide evolutionary rate was made with the Beast v1.10.4 program
159 package (Suchard et al. 2018). Analyses were run at the CIPRES Science Gateway server
160 (Miller et al. 2010). Three hundred and twelve sequences without indeterminations
161 corresponding to the nsp3 (5835nt) and S (3822nt) genes were randomly selected from
162 dataset 1. The sequences represent all the fortnights and most of the geographical locations
163 sampled until April 17. Temporal calibration was performed by date of sampling. The
164 appropriate evolutionary model was selected as described above for phylogenetic signal
165 analysis. The TIM model of nucleotide substitution was used for nsp3 and, the HKY model of
166 nucleotide substitution for S. The analysis was carried out under a relaxed (uncorrelated
167 lognormal) molecular clock model suggest by Duchene & col. (Duchene et al. 2020) and with
168 an exponential demographic, proper for early viral samples from an outbreak (Grassly &
169 Fraser, 2008). Independent runs were performed for each dataset and a Markov chain Monte
170 Carlo (MCMC) with a length of 1.3×10^9 steps, sampling every 1.3×10^6 steps, was set. The
171 convergence of the "meanRate" parameters [effective sample size (ESS) ≥ 200 , burn-in 10%]
172 was verified with Tracer v1.7.1 (Rambaut et al. 2018). Additionally, in order to verify the
173 obtained results, 15 independent replicates of the analysis were performed with the time
174 calibration information (date of sampling) randomized as described by Rieux & Khatchikian,
175 2017 (Rieux & Khatchikian, 2017). Finally, the obtained parameters for real data and the
176 randomized replicates were compared.

177

178 **3. Results**

179 *3.1 Phylogenetic signal*

180 Using bioinformatics tools, a phylogenetic signal study was carried out in order to identify the
181 most informative SARS-CoV-2 genomic regions. The likelihood mapping analysis showed that
182 most genes has very poor phylogenetic signal with high values in central region which
183 represents the area of unresolved quartets (Figure 1). Accordingly, genes could be separated
184 into three groups. A group with little or no phylogenetic signal (E, Orf6, Orf8, nsp1, and nsp14),
185 a second group with low phylogenetic signal (Orf3a and N), and a last group with relatively
186 more phylogenetic signal (S and nsp3) but still low to be considered a robust one (unresolved
187 quartets >40%).

188

189 *3.2 Analysis of amino acid substitutions*

190 The analysis of amino acids substitutions by fortnights was useful to study the viral
191 evolutionary dynamics in the context of the beginning of the pandemic. By analyzing different
192 time periods amino acid sequences, changes were observed in 5 out of 9 regions and only in
193 14 out of the 4975 (0.28%) evaluated residues. In most of the regions, except nsp1, nsp14, E,
194 and Orf6, 2 to 6 amino acids were selected since FN3 and remain unchanged until the end of
195 the follow up period (Table 2). Particularly, in Orf8 region, early selection of two amino acid
196 substitutions (V62L and S84L) was observed from FN2. On the other hand, in the S region,
197 the D614G substitution started with less than 2% in FN3 and FN4 and reached 88% in the last
198 fortnight. In a similar way, the Q57H (Orf3a) substitution went from 6% to 34% while S84L
199 start to be selected in FN2 and reached 94% by FN8. The R203K and G204R substitutions of
200 the N region was selected in FN4 and increased their population proportion with values greater
201 than 20% towards the end of the follow up period. Moreover, selection of a great number of
202 sporadic substitutions remaining in the population for a short period (1-3 fortnights) was

203 observed in the nine analyzed regions. Indeed, 333 (6.83%) of the analyzed positions
204 presented at least one substitution throughout the eight fortnights. Table 3 summarizes the
205 number of variable positions, number of mutations, and number of sequences with mutations
206 by region.

207

208 *3.3 Bayesian coalescence analysis*

209 In this study, trees were performed by Bayesian analysis instead of by distance, likelihood, or
210 parsimony methods. Consistently with the phylogenetic signal analysis, trees for nsp1, E, and
211 Orf6 showed a star-like topology. Nevertheless, different proportions of clades formation could
212 be observed in trees of Orf8, nsp14, Orf3a, N, S, and nsp3 regions (Figure 2). Finally, from
213 mentioned regions, nsp3 and S showed a better clade constitution. This analysis allowed to
214 differentiate regions presenting a diversification process (nsp3, nsp14, Orf3a, S, Orf8, and N)
215 from those that even after four months showed an incipient one (nsp1, E, and Orf6).
216 Furthermore, this nucleotide analysis is complemented by the previous study of amino acid
217 variations in each region. However, it is important to note that due to the low phylogenetic
218 signal observed for each region, results can only be considered as preliminary.

219

220 *3.4 Evolutionary rate*

221 Nsp3 and S sequences were selected to perform the evolutionary rate analysis since both
222 regions provided the best phylogenetic information among studied regions. The observed
223 evolutionary rate for nsp3 protein of SARS-CoV-2 was estimated to be 1.37×10^{-3} (ESS 782)
224 nucleotide substitutions per site per year (s/s/y) (95% HPD interval 9.16×10^{-4} to 1.91×10^{-3}).
225 On the other hand, the corresponding figures for S were estimated to be 2.19×10^{-3} (ESS 383)
226 nucleotide s/s/y (95% HPD interval 3.19×10^{-3} to 1.29×10^{-3}). In both genomic regions, date-
227 randomization analyses showed no overlapping between the 95% HPD substitution-rate

228 intervals obtained from real data and from date-randomized datasets. This fact suggests that
229 the original dataset has enough temporal signal to perform analyses with temporal calibration
230 based on tip-dates (Figure 3).

231 **4. Discussion**

232 The phylogenetic characterization of an emerging virus is crucial to understand the way the
233 virus and the pandemic will evolve. Thereby, a detailed study of the SARS CoV-2 genome
234 allows, on the one hand, to contribute to the knowledge of viral diversity in order to detect the
235 most suitable regions to be used as antivirals or vaccines targets. On the other hand, the large
236 amount of information that is continuously generated, is allowing studying the SARS CoV-2
237 genome and describing a new viral real time evolution like never before.

238 In the present study, the molecular evolution and viral lineages of SARS-CoV-2 in nine
239 genomic regions, during eight successive fortnights, was analyzed using three different
240 approaches: phylogenetic signal assessment, emergence of amino acid substitutions, and
241 Bayesian evolutionary rate estimation. In this context, the observed phylogenetic signals of
242 nine coding regions were very low and the obtained trees were consistent with this finding,
243 showing star-like topologies in some viral regions (nsp1, E, and Orf6). However, after a four
244 months evolution period, it was possible to identify regions (nsp3, S, Orf3a, Orf8, and N)
245 revealing an incipient formation of viral lineages, despite the phylogenetic signal, both at the
246 nucleotide and amino acid levels from FN3. Based on these findings, the SARS-CoV-2
247 evolutionary rate was estimated, for the first time, for the two regions showing higher variability
248 (S and nsp3).

249 As regards the phylogenetic signal, several simulation studies has proven that for a set of
250 sequences to be considered robust, the central and lateral areas representing the unresolved
251 quartets, must not be greater than 40% (Strimmer & von Haeseler, 1997). In this regard, none
252 of the nine analyzed regions met this requirement. Three regions (E, nsp1, and Orf6)
253 presented values of 100% unresolved quartets. Most regions (nsp14, Orf3a, Orf8, and N)
254 reached values higher than 85%. Only in regions nsp3 and S the number of unresolved
255 quartets dropped to ~ 60%. Thus, despite being a virus with an RNA genome, the short time

256 elapsed since its emergence, and possibly genetic restrictions have led to a constrained
257 evolution of SARS-CoV-2 in these months. For this reason, it is expected that trees generated
258 from SARS-CoV-2 partial sequences in the first months of the pandemic are unreliable for
259 defining clades. Therefore, they should be analyzed with great caution.

260 Since Bayesian analysis allows to infer phylogenetic patterns from tree distributions, it
261 represents a more reliable tool to compare different evolutionary behaviors. Bayesian analysis
262 helps to obtain a tree topology that is closer to reality in the current conditions of SARS-CoV-
263 2 pandemic (Drummond et al. 2006). The phylogenetic analysis for nsp1, E, and Orf6 regions
264 confirmed the star-like topologies in accordance to a lower diversification of these regions
265 using the sequences available up to FN8 (Figure 2). Trees generated from nsp14 and Orf8
266 are at an intermediate point, where the formation of small clusters can be observed. In fact, a
267 mutation at position 28,144 (Orf8: C251T, S84L) has been proposed as a possible marker for
268 viral classification (Tang et al. 2020, Yin, 2020). Finally, trees obtained from regions Orf3a, N,
269 nsp3, and S showed the best clade formation. Indeed, in the most variable regions nsp3 and
270 S, it can be clearly seen that sequences are separated into two large groups. Despite the
271 aforementioned for the nsp3 and S regions, even clusters with very high support values should
272 be taken with precaution and longer periods should be considered to obtain more accurate
273 phylogeny data. However, even when data are not the most accurate to study the spread or
274 clade formation (Mavian et al. 2020, Sanchez-Pacheco et al. 2020), they provide a good
275 representation of the way the virus is evolving.

276 The analysis of amino acids frequencies allowed identifying different degree of region
277 conservation throughout the viral genome as a consequence of positive and negative
278 pressures. In particular, nsp3, S, Orf8, and N showed some substitutions in high frequencies.
279 This would indicate, as other authors previously report, the frequent circulation of
280 polymorphisms due to significant positive pressure (Cagliani et al. 2020, Issa et al. 2020, Tang

281 et al. 2020). Additionally, since S and N are among candidates to be used in the formulation
282 of vaccines and antibody treatment, it will be important to monitor these substitutions in
283 different geographic regions in order to improve treatment and vaccination efficacy (Ahmed et
284 al. 2020, Callaway, 2020, Koyama et al. 2020). Contrarily, in regions nsp1, nsp14, E, and Orf6,
285 no substitutions were selected and lasted during the first 4 months of the pandemic. This would
286 suggest that these are regions with constraints to change due to the great negative selection
287 pressure, as it has been recently reported (Cagliani et al. 2020).

288 In the present study, the evolutionary rate for SARS-CoV-2 genes was estimated by analyzing
289 a large number of sequences, which were carefully curated and had a good temporal and
290 spatial structure. Additionally, the most phylogenetically informative regions of the genome
291 (nsp3 and S) were used for analysis, reinforcing the results confidence. Previous studies on
292 SARS-CoV-2 have reported similar data ranging from 1.79×10^{-3} to 6.58×10^{-3} s/s/y for the
293 complete genome (Giovanetti et al. 2020, Li et al. 2020a). However, in both articles, small
294 datasets of complete genomes were used (N=32 and 54, respectively). As studies were
295 performed early in the outbreak and due to datasets temporal structure, analysis could have
296 led to less precise estimates of the evolutionary rate (Duchene et al. 2020). On the other hand,
297 another study analyzing 7,666 sequences has obtained different results with a remarkably low
298 evolutionary rate (6×10^{-4} nucleotide/genome/year) (van Dorp et al. 2020). Additionally, tests
299 randomizing the dates of nsp3 and S datasets were carried out; they showed that these partial
300 genomic regions have enough temporal signal. In this context, our results (range from 1.37 to
301 2.19×10^{-3} s/s/y) are in close agreement with those published for SARS-CoV genome, which
302 have been estimated between 0.80 to 3.01×10^{-3} s/s/y (The Chinese SARS Molecular
303 Epidemiology Consortium, 2004, Vega et al. 2004, Zhao et al. 2004). Moreover, our values
304 are in the same magnitude order as other RNA viruses (Sanjuán, 2012). Even though we

305 should be cautious with these results interpretation, the date-randomization analysis indicated
306 a robust temporal signal.

307 Despite limitations of the evolutionary study of an emerging virus, where the selection
308 pressures are still low and therefore its variability is also low, this work has a great strength: it
309 lies on the extremely careful selection of a big sequence dataset to be analyze. First, it was
310 considered selected sequences to have a good temporal signal and spatial (geographic)
311 structure. Secondly, much attention was paid to the elimination of sequences with low
312 coverage and indeterminacies that could generate a noise for the phylogenetic analysis of a
313 virus that is beginning to evolve in a new host.

314 The appearance of a new virus means an adaptation challenge. The SARS-CoV-2 overcome
315 the spill stage and shows a significantly higher spread than SARS-CoV and MERS-CoV, thus
316 becoming itself the most important pandemic of the century. In this context, the results
317 obtained in this work about the variable diversity of nine crucial viral regions and the
318 determination of the evolutionary rate, are consequently decisive to understanding essential
319 feature of viral emergence. Nevertheless, monitoring SARS-CoV-2 population will be required
320 to determine the evolutionary course of new mutations as well as to understand the way they
321 affect viral fitness in human hosts.

322

323

324

325

326

327

328

329

330 **5. Acknowledgements**

331 MJP, LM, DMF, and FAD are members of the National Research Council (CONICET). We
332 would like to thank to the researchers who generated and shared the sequencing data from
333 GISAID (<https://www.gisaid.org/>) and Mrs. Silvina Heisecke from CEMIC-CONICET for
334 providing language assistance.

335

336 **6. REFERENCES**

- 337 [1] Ahmed S.F., Quadeer A.A. & McKay M.R (2020). Preliminary Identification of Potential
338 Vaccine Targets for the COVID-19 Coronavirus (SARS-CoV-2) Based on SARS-CoV
339 Immunological Studies. *Viruses*, 12, 254. <https://doi.org/10.3390/v12030254>
- 340 [2] Benvenuto D., Giovanetti M., Salemi M., Prosperi M., De Flora C., Junior Alcantara L.C.,
341 et al (2020). The global spread of 2019-nCoV: a molecular evolutionary analysis. *Pathogens*
342 *and Global Health*, 114, 64-67. <https://doi.org/10.1080/20477724.2020.1725339>
- 343 [3] Cagliani R., Forni D., Clerici M. & Sironi M (2020). Computational inference of selection
344 underlying the evolution of the novel coronavirus, SARS-CoV-2. *Journal of Virology*, [Preprint].
345 <https://doi.org/10.1128/JVI.00411-20>
- 346 [4] Callaway E (2020). The race for coronavirus vaccines: a graphical guide. *Nature*, 580, 576-
347 577. <https://doi.org/10.1038/d41586-020-01221-y>
- 348 [5] Cavanagh D. Coronavirus avian infectious bronchitis virus (2007). *Veterinary Research*,
349 38, 281-297. <https://doi.org/10.1051/vetres:2006055>
- 350 [6] Coronaviridae Study Group of the International Committee on Taxonomy of Viruses (2020).
351 The species *Severe acute respiratory syndrome-related coronavirus*: classifying 2019-nCoV
352 and naming it SARS-CoV-2. *Nature Microbiology*, 5, 536-544. [https://doi.org/10.1038/s41564-](https://doi.org/10.1038/s41564-020-0695-z)
353 [020-0695-z](https://doi.org/10.1038/s41564-020-0695-z)
- 354 [7] Cui J., Li F. & Shi Z.L. (2019). Origin and evolution of pathogenic coronaviruses. *Nature*
355 *Reviews Microbiology*, 17, 181-192. <https://doi.org/10.1038/s41579-018-0118-9>
- 356 [8] Darriba D., Taboada G.L., Doallo R. & Posada D (2012). jModelTest 2: more models, new
357 heuristics and parallel computing. *Nature Methods*, 9, 772. <https://doi.org/10.1038/nmeth.2109>
- 358 [9] Drummond A.J., Ho S.Y., Phillips M.J. & Rambaut A (2006). Relaxed phylogenetics and
359 dating with confidence. *PLoS Biology*, 4, e88. <https://doi.org/10.1371/journal.pbio.0040088>

- 360 [10] Duchene S., Featherstone L., Haritopoulou-Sinanidou M., Rambaut A., Lemey P., & Baele
361 G (2020). Temporal signal and the phylodynamic threshold of SARS-CoV-2. *bioRxiv*,
362 [Preprint]. <https://doi.org/10.1101/2020.05.04.077735>
- 363 [12] Giovanetti M., Benvenuto D., Angeletti S. & Ciccozzi M (2020). The first two cases of
364 2019-nCoV in Italy: Where they come from? *Journal of Medical Virology*, 92, 518-521.
365 <https://doi.org/10.1002/jmv.25699>
- 366 [13] Grassly NC & Fraser C. (2008). Mathematical models of infectious disease transmission.
367 *Nat Rev Microbiol.* 6, 477-487. <https://doi.org/10.1038/nrmicro1845>
- 368 [14] Ismail M.M., Tang A.Y. & Saif Y.M. (2003). Pathogenicity of turkey coronavirus in turkeys
369 and chickens. *Avian Diseases*, 47, 515-522. <https://doi.org/10.1637/5917>
- 370 [15] Issa E., Merhi G., Panossian B., Salloum T. & Tokajian S (2020). S.SARS-CoV-2 and
371 ORF3a: Nonsynonymous Mutations, Functional Domains, and Viral Pathogenesis. *mSystems*,
372 [Preprint]. <https://doi.org/10.1128/mSystems.00266-20>
- 373 [16] Koyama T., Weeraratne D., Snowdon J.L. & Parida L (2020). Emergence of Drift Variants
374 That May Affect COVID-19 Vaccine Development and Antibody Treatment.
375 *Pathogens*, 9, 324. <https://doi.org/10.20944/preprints202004.0024.v1>
- 376 [17] Li X., Wang W., Zhao X., Zai J., Zhao Q., Li Y., et al (2020a). Transmission dynamics and
377 evolutionary history of 2019-nCoV. *Journal of Medical Virology*, 92, 501-511.
378 <https://doi.org/10.1002/jmv.25701>
- 379 [18] Li Q., Guan X., Wu P., Wang X., Zhou L., Tong Y., et al (2020b). Early Transmission
380 Dynamics in Wuhan, China, of Novel Coronavirus-Infected Pneumonia. *The New England*
381 *Journal of Medicine*, 382, 1199-1207. <https://doi.org/10.1056/NEJMoa2001316>.
- 382 [19] Luk H.K.H., Li X., Fung J., Lau S.K.P. & Woo P.C.Y (2019). Molecular epidemiology,
383 evolution and phylogeny of SARS coronavirus. *Infection Genetics and Evolution*, 71, 21-30.
384 <https://doi.org/10.1016/j.meegid.2019.03.001>

- 385 [20] Mavian C., Marini S., Prosperi M. & Salemi M (2020). A snapshot of SARS-CoV-2 genome
386 availability up to 30th March, 2020 and its implications. *bioRxiv*, [Preprint].
387 <https://doi.org/10.1101/2020.04.01.020594>
- 388 [21] Miller M.A., Pfeiffer W. & Schwartz T (2010). Creating the CIPRES Science Gateway for
389 inference of large phylogenetic trees. *Gateway Computing Environments Workshop*, 1-8.
390 <https://doi.org/10.1109/GCE.2010.5676129>
- 391 [22] Phan T. (2020). Genetic diversity and evolution of SARS-CoV-2. *Infection Genetics and*
392 *Evolution*, 81, 104260. <https://doi.org/10.1016/j.meegid.2020.104260>
- 393 [23] Rambaut A., Drummond A.J., Xie D., Baele G. & Suchard M.A. (2018). Posterior
394 summarization in Bayesian phylogenetics using Tracer 1.7. *Systematic Biology*, 67, 901-904.
395 <https://doi.org/10.1093/sysbio/syy032>
- 396 [24] Rieux A. & Khatchikian C.E. (2017). tipdatingbeast: an r package to assist the
397 implementation of phylogenetic tip-dating tests using beast. *Molecular Ecology Resources*,
398 17, 608-613. <https://doi.org/10.1111/1755-0998.12603>
- 399 [25] Ronquist F., Teslenko M., van der Mark P., Ayres D.L., Darling A., Höhna S., et al. (2012).
400 MrBayes 3.2: efficient Bayesian phylogenetic inference and model choice across a large
401 model space. *Systematic Biology*, 61, 539-542. <https://doi.org/10.1093/sysbio/sys029>
- 402 [26] Sánchez-Pacheco S.J., Kong S., Pulido-Santacruz P., Murphy R.W. & Kubatko L. (2020).
403 Median-joining network analysis of SARS-CoV-2 genomes is neither phylogenetic nor
404 evolutionary. *Proceedings of the National Academy of Sciences of the USA*, 117, 9241–9243.
405 <https://doi.org/10.1073/pnas.2007062117>
- 406 [27] Sanjuán R (2012). From molecular genetics to phylodynamics: evolutionary relevance of
407 mutation rates across viruses. *PLoS Pathogens*, 8, e1002685. [https://doi.org/](https://doi.org/10.1371/journal.ppat.1002685)
408 [10.1371/journal.ppat.1002685](https://doi.org/10.1371/journal.ppat.1002685)

- 409 [28] Schmidt H.A., Strimmer K., Vingron M. & Von Haeseler A (2002). TREE-PUZZLE:
410 Maximum likelihood phylogenetic analysis using quartets and parallel computing.
411 *Bioinformatics*, 18, 502-504. <https://doi.org/10.1093/bioinformatics/18.3.502>
- 412 [29] Strimmer K. & von Haeseler A (1997). Likelihood-mapping: A simple method to visualize
413 phylogenetic content of a sequence alignment. *Proceedings of the National Academy of*
414 *Sciences of the USA*, 94, 6815-6819. <https://doi.org/10.1073/pnas.94.13.6815>
- 415 [30] Su S., Wong G., Shi W., Liu J., Lai A.C.K., Zhou J., et al (2016). Epidemiology, genetic
416 recombination, and pathogenesis of coronaviruses. *Trends in Microbiology*, 24, 490-502.
417 <https://doi.org/10.1016/j.tim.2016.03.003>
- 418 [31] Suchard M.A., Lemey P., Baele G., Ayres D.L., Drummond A.J. & Rambaut A (2018).
419 Bayesian phylogenetic and phylodynamic data integration using BEAST 1.10. *Virus Evolution*,
420 4, vey016. <https://doi.org/10.1093/ve/vey016>
- 421 [32] Tang X., Wu C., Li X., Song Y., Yao X., Wu X., et al (2020). On the origin and continuing
422 evolution of SARS-CoV-2. *National Science Review*, 0, 1-12.
423 <https://doi.org/10.1093/nsr/nwaa036>
- 424 [33] The Chinese SARS Molecular Epidemiology Consortium (2004). Molecular Evolution of
425 the SARS Coronavirus During the Course of the SARS Epidemic in China. *Science*, 303,
426 1666-1669. <https://doi.org/10.1126/science.1092002>
- 427 [34] van Dorp L., Acman M., Richard D., Shaw L.P., Ford C.E., Ormond L., et al (2020).
428 Emergence of genomic diversity and recurrent mutations in SARS-CoV-2. *Infection Genetics*
429 *and Evolution*, 5, 104351. <https://doi.org/10.1016/j.meegid.2020.104351>
- 430 [35] Vega V.B., Ruan Y., Liu J., Lee W.H., Wei C.L., Se-Thoe S.Y., et al (2004). Mutational
431 dynamics of the SARS coronavirus in cell culture and human populations isolated in 2003.
432 *BMC Infectious Diseases*, 4, 32. <https://doi.org/10.1186/1471-2334-4-32>

- 433 [36] World Health Organization, 2020. Coronavirus disease (COVID-19) Situation Report – 118.
434 Retrieved from: [https://www.who.int/docs/default-source/coronaviruse/situation-](https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200517-covid-19-sitrep-118.pdf?sfvrsn=21c0d4fe_6)
435 [reports/20200517-covid-19-sitrep-118.pdf?sfvrsn=21c0d4fe_6](https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200517-covid-19-sitrep-118.pdf?sfvrsn=21c0d4fe_6) (17 July 2020, date last
436 accessed).
- 437 [37] Yin C (2020). Genotyping coronavirus SARS-CoV-2: methods and implications.
438 *Genomics*, 30318-30319, [Preprint]. <https://doi.org/10.1016/j.ygeno.2020.04.016>
- 439 [38] Zhao Z., Li H., Wu X., Zhong Y., Zhang K., Zhang Y.P., et al (2004). Moderate mutation
440 rate in the SARS coronavirus genome and its implications. *BMC Evolutionary Biology*, 4, 21.
441 <https://doi.org/10.1186/1471-2148-4-21>
- 442 [39] Zhou P., Yang X.L., Wang X.G., Hu B., Zhang L., Zhang W. et al (2020). A pneumonia
443 outbreak associated with a new coronavirus of probable bat origin. *Nature*, 579, 270-273.
444 <https://doi.org/10.1038/s41586-020-2012-7>
- 445 [40] Zhu N., Zhang D., Wang W., Li X., Yang B., Song J., et al. (2020). A Novel Coronavirus
446 from Patients with Pneumonia in China, 2019. *The New England Journal of Medicine*, 382,
447 727-733. <https://doi.org/10.1056/NEJMoa2001017>
- 448
- 449
- 450
- 451
- 452
- 453
- 454
- 455
- 456
- 457

458 **Table 1.** Number of SARS-CoV-2 sequences by fortnight (Temporal structure)

Fortnight	Date	Median of analyzed sequences (Q1-Q3)
FN1	12/24/2019 to 12/31/2019	15
FN2	01/01/2020 to 01/15/2020	19
FN3	01/16/2020 to 01/31/2020	145 (136-145.5)
FN4	02/01/2020 to 02/15/2020	119 (113-120)
FN5	02/16/2020 to 03/02/2020	258 (247-259)
FN6	03/03/2020 to 03/17/2020	403 (390-406)
FN7	03/18/2020 to 04/01/2020	447 (416-450)
FN8	04/02/2020 to 04/17/2020	199 (197-201)
TOTAL		1488 to 1616

459 FN: Fortnight; Q1=quartile 1, Q3=quartile 3. The total number of sequences is variable depending on
460 the analyzed region (nsp1, 1608; nsp3, 1511; nsp14, 1550; S, 1488; Orf3a, 1600; E, 1615; Orf6, 1616;
461 Orf8, 1612; and N, 1610)
462

463 **Table 2.** Amino acids selected by region and fortnight. The number indicates the amino acids
 464 location in its protein.

Region	Amino acid substitution	Amino acid percentage by FN							
		FN1	FN2	FN3	FN4	FN5	FN6	FN7	FN8
nsp3	A58T	0	0	0	1.0	6.0	3.0	3.0	2.5
	P135L	0	0	0.8	0	0	1.5	0.5	2.5
S	D614G	0	0	1.5	1.8	37.0	64.0	75.0	88.0
Orf3a	Q75H	0	0	0	0	6.0	22.0	23.0	34.0
	G196V	0	0	0	0	0.8	4.0	0.9	0.5
	G251V	0	0	8.0	24.0	8.0	9.0	10.0	3.0
Orf8	V62L	0	5.0	1.0	3.3	0.0	1.5	1.3	3.0
	S84L	0	58.0	63.0	79.0	79.0	82.0	93.0	94.0
N	P13L	0	0	0	0	1.0	1.0	2.5	0.5
	S197L	0	0	0	0	1.1	5.0	0.9	0.5
	S202N	0	0	3.5	4.2	0	0.5	2.2	2.5
	R203K	0	0	0	0	17.0	19.0	24.0	23.0
	G204R	0	0	0	0	17.0	19.0	24.0	23.0
	I292T	0	0	0	0	2.0	0.2	0.2	0.5

465 Only regions where amino acid change was selected and remained until the last analyzed fortnight are
 466 shown. FN: Fortnight; aa: amino acid
 467
 468

469 **Table 3.** Number of variable positions, number of mutations, and number of sequences with
470 mutation by region

Region	N° of variable aa positions (%)	N° of aa substitutions	N° of sequences with aa substitutions (%)
nsp1 (180aa)	3 (1.7)	37	37 (2.4)
nsp3 (1945aa)	158 (8.1)	322	294 (19.3)
nsp14 (527aa)	6 (1.4)	83	83 (5.5)
S (1273aa)	76 (5.9)	1013	904 (59.4)
Orf3a (275aa)	11 (4)	491	468 (30.7)
E (75aa)	5 (6.7)	6	6 (0.4)
Orf6 (60aa)	7 (11.6)	9	9 (0.6)
Orf8 (121aa)	14 (11.6)	312	288 (18.9)
N (419aa)	53 (12.6)	760	470 (30.9)
Total (4875aa)	333 (6.8)	3033	-

471 aa: amino acid

472

473

474

475

476

477

478

479

480

481

482

483

484 **Figure legends**

485

486 **Figure 1:** Phylogenetic signal for SARS-CoV-2 datasets. Presence of phylogenetic signal was
487 evaluated by likelihood mapping, unresolved quartets (center) and partly resolved quartets
488 (edges) for genomes available on April 17 for the nine analyzed regions: nsp1 (29 sequences),
489 nsp3 (225 sequences), nsp14 (65 sequences), S (183 sequences), Orf3a (74 sequences), E
490 (11 sequences), Orf6 (12 sequences), Orf8 (23 sequences), and N (113 sequences).
491 Presence of strong phylogenetic signal (<40% unresolved quartets) was not reached for any
492 region.

493

494 **Figure 2:** Bayesian trees of 29 sequences of nsp1 (540nt), 225 sequences of nsp3 (5835nt),
495 65 sequences of nsp14 (1581nt), 183 sequences of S (3822nt), 74 sequences of Orf3a
496 (828nt), 11 sequences of E (228nt), 12 sequences of Orf6 (186nt), 23 sequences of Orf8
497 (366nt), and 113 sequences of N (1260nt). Scale bar represents substitutions per site.

498

499 **Figure 3:** Comparison of the evolutionary rates estimated using BEAST for the original dataset
500 and the date-randomized datasets (312 sequences). This analysis was performed for regions
501 nsp3 (5835nt) and S (3822nt). s.s.y = substitutions/site/year.

502

503

504

505

506

507

508

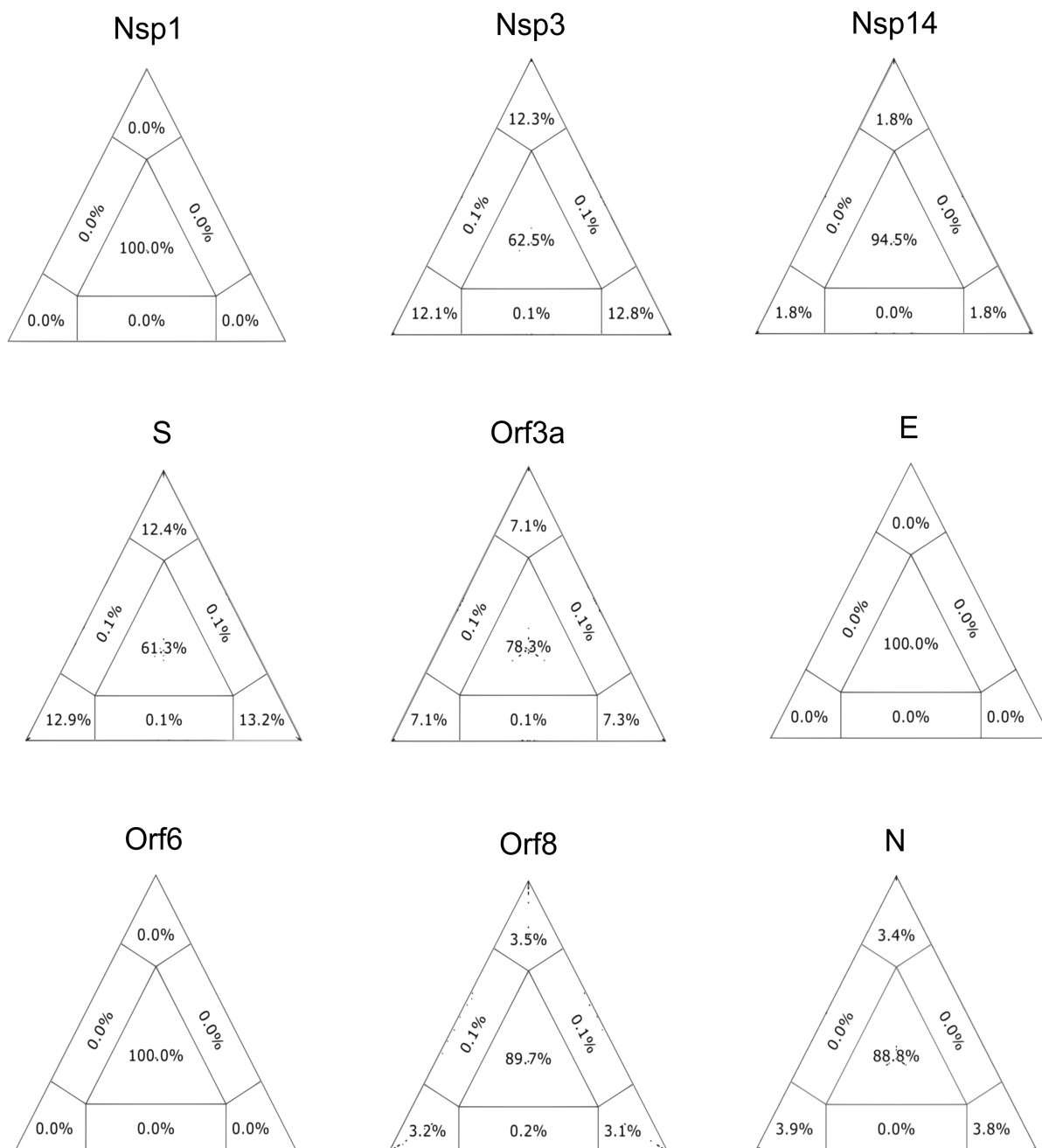


Figure 1

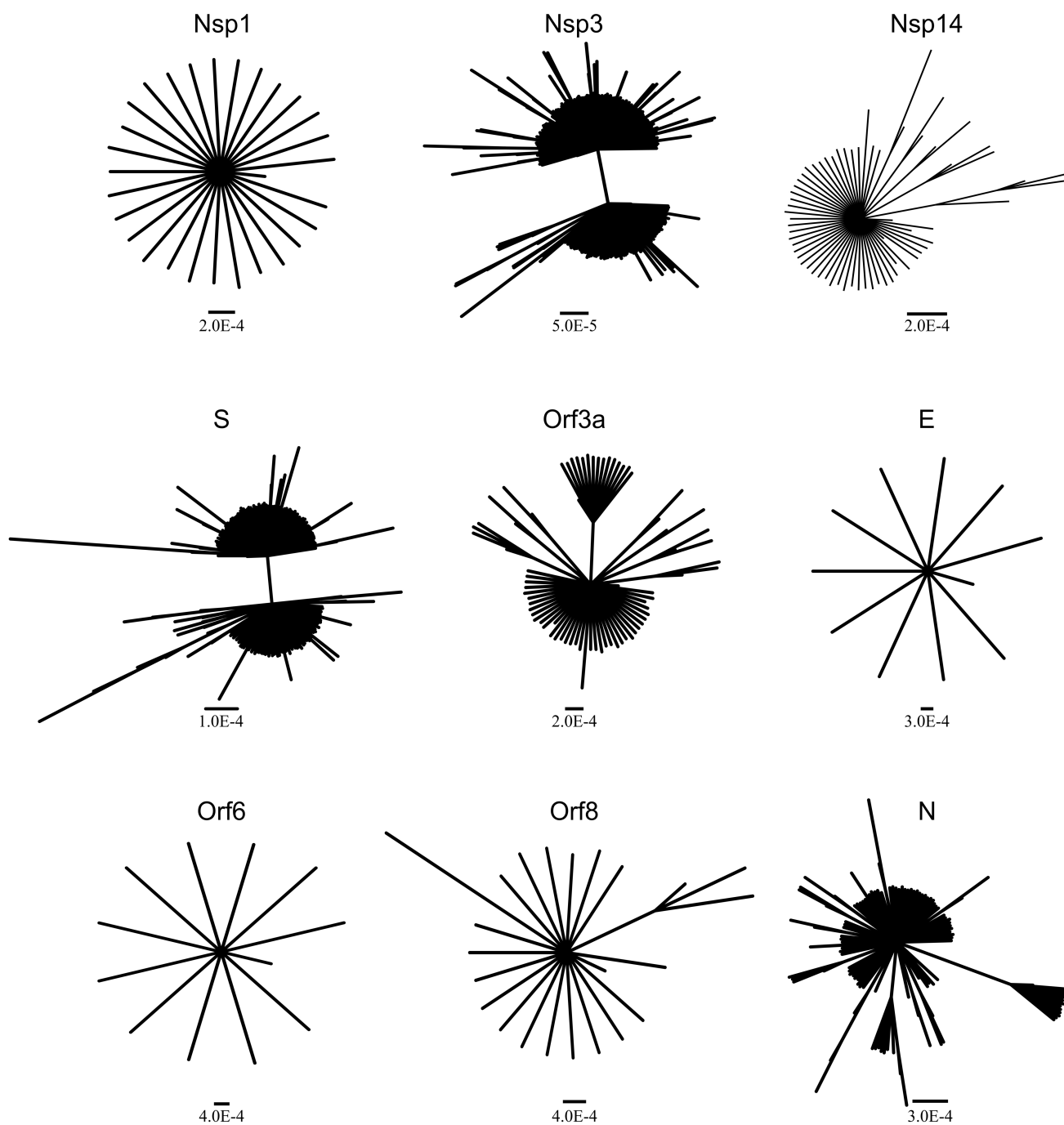


Figure 2

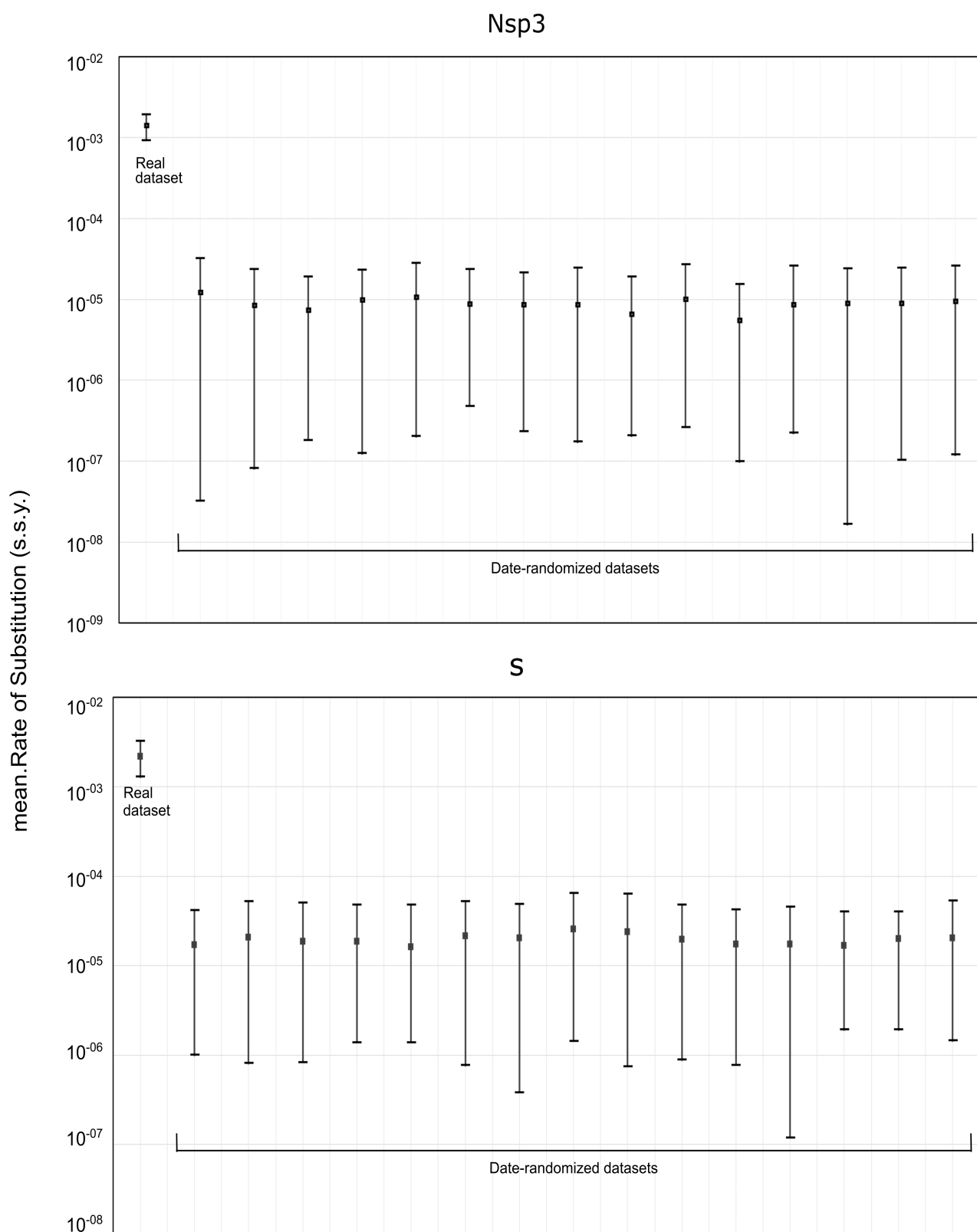


Figure 3