# The origin of a new human virus: phylogenetic analysis of the evolution of sars-cov-2 — Source link ↗

Matías J. Pereson, Matías J. Pereson, Laura Noelia Mojsiejczuk, Laura Noelia Mojsiejczuk ...+5 more authors

**Institutions:** National Scientific and Technical Research Council, University of Buenos Aires

Related papers:

- Phylogenetic Analysis Of SARS-CoV-2 In The First Months Since Its Emergence

- Phylogenetic analysis of SARS-CoV-2 in the first few months since its emergence.

- Emerging phylogenetic structure of the SARS-CoV-2 pandemic

- Comparative genomics suggests limited variability and similar evolutionary patterns between major clades of SARS-Cov-2

- Genomic and evolutionary comparison between SARS-CoV-2 and other human coronaviruses. (Special issue on Covid-19.)

Share this paper: ❍ 🐦 in ✉

1  **TITLE PAGE**

2  **Title:** PHYLOGENETIC ANALYSIS OF SARS-COV-2 IN THE FIRST MONTHS SINCE ITS

3  EMERGENCE

4

5  **Authors:** Matías J. PERESON[a,b], Laura MOJSIEJCZUK[a,b], Alfredo P. MARTÍNEZ[c], Diego M.

6  FLICHMAN[b,d], Gabriel H. GARCIA[a], Federico A. DI LELLO[a,b]

7

8  **Affiliations:**

9  [a]Universidad de Buenos Aires. Facultad de Farmacia y Bioquímica. Instituto de

10  Investigaciones en Bacteriología y Virología Molecular (IBaViM). Buenos Aires, Argentina.

11  [b]Consejo Nacional de Investigaciones Científicas y Técnicas (CONICET), Ciudad Autónoma

12  de Buenos Aires, Argentina.

13  [c]Virology Section, Centro de Educación Médica e Investigaciones Clínicas Norberto Quirno

14  "CEMIC". Buenos Aires, Argentina.

15  [d]Instituto de Investigaciones Biomédicas en Retrovirus y Síndrome de Inmunodeficiencia

16  Adquirida (INBIRS) – Consejo Nacional de Investigaciones Científicas y Técnicas

17  (CONICET), Universidad de Buenos Aires, Buenos Aires, Argentina.

18

19  **Corresponding author:**

20  Dr. Federico Alejandro Di Lello, Facultad de Farmacia y Bioquímica, Universidad de Buenos

21  Aires, Instituto de Investigaciones en Bacteriología y Virología Molecular (IBaViM). Junín

22  956, 4º piso, (1113), Ciudad Autónoma de Buenos Aires, Argentina.

23  Phone: +54 11 5287 4472, Fax: +54 11 5287 4662, E-mail: fadilello@ffyb.uba.ar

24

25  **Running Title**: Phylogenetic analysis and evolution of SARS-CoV-2

26  **ABSTRACT**

27  During the first months of SARS-CoV-2 evolution in a new host, contrasting hypotheses

28  have been proposed about the way the virus has evolved and diversified worldwide. The aim

29  of this study was to perform a comprehensive evolutionary analysis to describe the human

30  outbreak and the evolutionary rate of different genomic regions of SARS-CoV-2.

31  The molecular evolution in nine genomic regions of SARS-CoV-2 was analyzed using three

32  different approaches: phylogenetic signal assessment, emergence of amino acid

33  substitutions, and Bayesian evolutionary rate estimation in eight successive fortnights since

34  the virus emergence.

35  All observed phylogenetic signals were very low and trees topologies were in agreement

36  with those signals. However, after four months of evolution, it was possible to identify

37  regions revealing an incipient viral lineages formation despite the low phylogenetic signal,

38  since fortnight 3. Finally, the SARS-CoV-2 evolutionary rate for regions nsp3 and S, the ones

39  presenting greater variability, was estimated to values of $1.37 \times 10^{-3}$ and $2.19 \times 10^{-3}$

40  substitution/site/year, respectively.

41  In conclusion, results obtained in this work about the variable diversity of crucial viral regions

42  and the determination of the evolutionary rate are consequently decisive to understand

43  essential feature of viral emergence. In turn, findings may allow characterizing for the first

44  time, the evolutionary rate of S protein that is crucial for vaccines development.

45

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

**Introduction**

47 

48 Coronaviruses belong to *Coronaviridae* family and have a single strand of positive-sense

49 RNA genome of 26 to 32 kb in length [1]. They have been identified in different avian hosts as

50 well as in various mammals including bats, mice, dogs, etc. [2,3]. Periodically, new

51 mammalian coronaviruses are identified. In late December 2019, Chinese health authorities

52 identified groups of patients with pneumonia of unknown cause in Wuhan, Hubei Province,

53 China [4]. The pathogen, a new coronavirus called SARS-CoV-2 [5], was identified by local

54 hospitals using a surveillance mechanism for "pneumonia of unknown etiology" [4,6,7]. The

55 pandemic has spread rapidly and, to date, more than 22 million confirmed cases and nearly

56 750,000 deaths have been reported in just over a six months period [8]. This rapid viral

57 spread raises interesting questions about the way its evolution is driven during the

58 pandemic. From the SARS-CoV-2 genome, 16 non-structural proteins (nsp1-16), 4 structural

59 proteins [spike (S), envelope (E), membrane (M) and nucleoprotein (N)], and other proteins

60 essential to complete the replication cycle are translated [9,10]. The large amount of

61 information currently available allows knowing, as never before, the real-time evolution

62 history of a virus since its interspecies jump [11]. Most studies published to date have

63 characterized the viral genome and evolution by analyzing complete genomes sequences

64 [12,13,14,15]. Despite this, until now, the viral genomic region providing the most accurate

65 information to characterize SARS-CoV-2, could not be established. This lack of information

66 prevent from investigating its molecular evolution and monitoring biological features affecting

67 the development of antiviral and vaccines. Therefore, the aim of this study was to perform a

68 comprehensive viral evolutionary analysis in order to describe the human outbreak and the

69 molecular evolution rate of different genomic regions of SARS-CoV-2.

70 **Materials and Methods**

71 *Datasets*

72 In order to generate a dataset representing different geographic regions and time evolution

73 of the SARS-CoV-2 pandemic from December 2019 to April 2020, data of all the complete

74 genome sequences available at GISAID (https: //www.gisaid.org /) on April 18, 2020 were

75 collected. Data inclusion criteria were: a.- complete genomes, b.- high coverage level, and

76 c.- human hosts only (no other animals, cell culture, or environmental samples). Complete

77 genomes were aligned using MAFFT against the Wuhan-Hu-1 reference genome

78 (NC_045512.2, EPI_ISL_402125). The resulting multiple sequence alignment (dataset 1)

79 was split in nine datasets corresponding to nine coding regions: a.- four structural proteins

80 [envelope (E), nucleocapsid (N), spike (S), Orf3a], b.- four nonstructural proteins (nsp1,

81 nsp3, Orf6, and nsp14), and c.- an unknown function protein (Orf8).

82 More than six thousand SARS-CoV-2 publicly available nucleotide sequences were

83 downloaded. After data selection according to the inclusion criteria, 1616 SARS-CoV-2

84 complete genomes were included in dataset 1. Sequences of this dataset 1 came from 55

85 countries belonging to the five continents as follow: Africa: 39 sequences, Americas: 383

86 sequences, Asia: 387 sequences, Europe: 686 sequences and Oceania: 121 sequences.

87 After elimination of sequences with indeterminate or ambiguous positions, the number of

88 analyzed sequences for each region were: nsp1, 1608; nsp3, 1511; nsp14, 1550; S, 1488;

89 Orf3a, 1600; E, 1615; Orf6, 1616; Orf8, 1612; and N, 1610. Finally, nucleotide sequences

90 were grouped by fortnight (FN) according to their collection date. Table 1 summarizes the

91 number of sequences per fortnight since the beginning of the pandemic up to FN 8. On the

92 other hand, Dataset 2 was created using only variable sequences of each region analyzed in

93 Dataset 1. Thus, Dataset 1 was used for the analysis of amino acid substitutions and

94    Dataset 2 was used for the phylogenetic signal analysis and the Bayesian coalescent trees

95    construction.

96

97    *Phylogenetic signal*

98    To determine the phylogenetic signal of each of the nine generated alignments, Likelihood

99    Mapping analyzes were carried out [16], using the Tree Puzzle v5.3 program [17] and the

100   Quartet puzzling algorithm. This algorithm allowed analyzing the tree topologies that can be

101   completely solved from all possible quartets of the n alignment sequences using maximum

102   likelihood. An alignment with defined tree values greater than 70-80% presents strong

103   support from the statistical point of view [17]. Identical sequences were also removed with

104   ElimDupes                                    (Available                                    at

105   https://www.hiv.lanl.gov/content/sequence/elimdupesv2/elimdupes.html) as they increase

106   computation time and provide no additional information about data phylogeny. The best-fit

107   evolutionary model to each dataset was selected based on the Bayesian Information

108   Criterion obtained with the JModelTest v2.1.10 software [18].

109

110   *Analysis of amino acid substitutions*

111   Entropy-One                                  (Available                                    at

112   https://www.hiv.lanl.gov/content/sequence/ENTROPY/entropy_one.html)    was    used    to

113   determining in dataset 1 the frequency of amino acids at each position for the nine genomic

114   regions analyzed and evaluating their permanence in the eight investigated fortnights.

115

116   *Bayesian coalescence and phylogenetic analysis*

117   To study the relationship between SARS-CoV-2 sequences, nine regions of the virus

118   genome were investigated by Bayesian analyses. Phylogenetic trees were constructed using

119 Bayesian inference with MrBayes v3.2.7a [19]. Each gene was analyzed independently with

120 the same dataset used for the phylogenetic signal analysis so that non-identical sequences

121 were included in the analysis. Analyses were run for five million generations and sampled

122 every 5000 generations. Convergence of parameters [effective sample size (ESS) ≥ 200,

123 with a 10% burn-in] was verified with Tracer v1.7.1 [20]. Phylogenetic trees were visualized

124 with FigTree v1.4.4.

125

126 *Evolutionary rate*

127 The estimation of the nucleotide evolutionary rate was made with the Beast v1.10.4 program

128 package [21]. Analyses were run at the CIPRES Science Gateway server [22]. Three hundred

129 and twelve sequences without indeterminations corresponding to the nsp3 (5835nt) and S

130 (3822nt) genes were randomly selected from dataset 1. The sequences represent all the

131 fortnights and most of the geographical locations sampled until April 17. Temporal calibration

132 was performed by date of sampling. The appropriate evolutionary model was selected as

133 described above for phylogenetic signal analysis. The TIM model of nucleotide substitution

134 was used for nsp3 and, the HKY model of nucleotide substitution for S. The analysis was

135 carried out under a relaxed (uncorrelated lognormal) molecular clock model suggest by

136 Duchene & col. [23] and with an exponential demographic, proper for early viral samples from

137 an outbreak [24]. Independent runs were performed for each dataset and a Markov chain

138 Monte Carlo (MCMC) with a length of $1.3 \times 10^9$ steps, sampling every $1.3 \times 10^6$ steps, was set.

139 The convergence of the "meanRate" parameters [effective sample size (ESS) ≥ 200, burn-in

140 10%] was verified with Tracer v1.7.1 [20]. Additionally, in order to verify the obtained results,

141 15 independent replicates of the analysis were performed with the time calibration

142 information (date of sampling) randomized as described by Rieux & Khatchikian, 2017 [25].

143 Finally, the obtained parameters for real data and the randomized replicates were compared.

144 **Results**

145 *Phylogenetic signal*

146 Using bioinformatics tools, a phylogenetic signal study was carried out in order to identify the

147 most informative SARS-CoV-2 genomic regions. The likelihood mapping analysis showed

148 that most genes has very poor phylogenetic signal with high values in central region which

149 represents the area of unresolved quartets (Figure 1). Accordingly, genes could be

150 separated into three groups. A group with little or no phylogenetic signal (E, Orf6, Orf8, nsp1,

151 and nsp14), a second group with low phylogenetic signal (Orf3a and N), and a last group

152 with relatively more phylogenetic signal (S and nsp3) but still low to be considered a robust

153 one (unresolved quartets >40%).

154

155 *Analysis of amino acid substitutions*

156 The analysis of amino acids substitutions by fortnights was useful to study the viral

157 evolutionary dynamics in the context of the beginning of the pandemic. By analyzing different

158 time periods amino acid sequences, changes were observed in 5 out of 9 genomic regions

159 and only in 14 out of the 4975 (0.28%) evaluated residues. In most of the regions, except

160 nsp1, nsp14, E, and Orf6, 2 to 6 amino acids were selected since FN3 and remain

161 unchanged until the end of the follow up period (Table 2). Particularly, in Orf8 region, early

162 selection of two amino acid substitutions (V62L and L84S) was observed from FN2. On the

163 other hand, in the S region, the D614G substitution started with less than 2% in FN3 and

164 FN4 and reached 88% in the last fortnight. In a similar way, the Q57H (Orf3a) substitution

165 went from 6% to 34% while L84S (Orf8) start to be selected in FN2 and reached 6% by FN8.

166 The R203K and G204R substitutions of the N region was selected in FN4 and increased

167 their population proportion with values greater than 20% towards the end of the follow up

168 period. Moreover, selection of a great number of sporadic substitutions remaining in the

169 population for a short period (1-3 fortnights) was observed in the nine analyzed regions.

170 Indeed, 333 (6.83%) of the analyzed positions presented at least one substitution throughout

171 the eight fortnights. Table 3 summarizes the number of variable positions, number of

172 mutations, and number of sequences with mutations by region.

173

174 *Bayesian coalescence analysis*

175 In this study, trees were performed by Bayesian analysis instead of by distance, likelihood,

176 or parsimony methods. Consistently with the phylogenetic signal analysis, trees for nsp1, E,

177 and Orf6 showed a star-like topology. Nevertheless, different proportions of clades formation

178 could be observed in trees of Orf8, nsp14, Orf3a, N, S, and nsp3 regions (Figure 2). Finally,

179 from mentioned regions, nsp3 and S showed a better clade constitution. This analysis

180 allowed to differentiate regions presenting a diversification process (nsp3, nsp14, Orf3a, S,

181 Orf8, and N) from those that even after four months showed an incipient one (nsp1, E, and

182 Orf6). Furthermore, this nucleotide analysis is complemented by the previous study of amino

183 acid variations in each region. However, it is important to note that due to the low

184 phylogenetic signal observed for each region, results can only be considered as preliminary.

185

186 *Evolutionary rate*

187 Nsp3 and S sequences were selected to perform the evolutionary rate analysis since both

188 regions provided the best phylogenetic information among studied regions. The observed

189 evolutionary rate for nsp3 protein of SARS-CoV-2 was estimated to be $1.37 \times 10^{-3}$ (ESS 782)

190 nucleotide substitutions per site per year (s/s/y) (95% HPD interval $9.16 \times 10^{-4}$ to $1.91 \times 10^{-3}$).

191 On the other hand, the corresponding figures for S were estimated to be $2.19 \times 10^{-3}$ (ESS

192 383) nucleotide s/s/y (95% HPD interval $3.19 \times 10^{-3}$ to $1.29 \times 10^{-3}$). In both genomic regions,

193 date-randomization analyses showed no overlapping between the 95% HPD substitution-

194    rate intervals obtained from real data and from date-randomized datasets. This fact suggests

195    that the original dataset has enough temporal signal to perform analyses with temporal

196    calibration based on tip-dates (Figure 3).

**Discussion**

197

198    The phylogenetic characterization of an emerging virus is crucial to understand the way the

199    virus and the pandemic will evolve. Thereby, a detailed study of the SARS CoV-2 genome

200    allows, on the one hand, to contribute to the knowledge of viral diversity in order to detect

201    the most suitable regions to be used as antivirals or vaccines targets. On the other hand, the

202    large amount of information that is continuously generated, is allowing studying the SARS

203    CoV-2 genome and describing a new viral real time evolution like never before.

204    In the present study, the molecular evolution and viral lineages of SARS-CoV-2 in nine

205    genomic regions, during eight successive fortnights, was analyzed using three different

206    approaches: phylogenetic signal assessment, emergence of amino acid substitutions, and

207    Bayesian evolutionary rate estimation. In this context, the observed phylogenetic signals of

208    nine coding regions were very low and the obtained trees were consistent with this finding,

209    showing star-like topologies in some viral regions (nsp1, E, and Orf6). However, after a four

210    months evolution period, it was possible to identify regions (nsp3, S, Orf3a, Orf8, and N)

211    revealing an incipient formation of viral lineages, despite the phylogenetic signal, both at the

212    nucleotide and amino acid levels from FN3. Based on these findings, the SARS-CoV-2

213    evolutionary rate was estimated, for the first time, for the two regions showing higher

214    variability (S and nsp3).

215    As regards the phylogenetic signal, several simulation studies has proven that for a set of

216    sequences to be considered robust, the central and lateral areas representing the

217    unresolved quartets, must not be greater than 40% [16]. In this regard, none of the nine

218    analyzed regions met this requirement. Three regions (E, nsp1, and Orf6) presented values

219    of 100% unresolved quartets. Most regions (nsp14, Orf3a, Orf8, and N) reached values

220    higher than 85%. Only in regions nsp3 and S the number of unresolved quartets dropped to

221    ~ 60%. Thus, despite being a virus with an RNA genome, the short time elapsed since its

222    emergence, and possibly genetic restrictions have led to a constrained evolution of SARS-

223    CoV-2 in these months. For this reason, it is expected that trees generated from SARS-CoV-

224    2 partial sequences in the first months of the pandemic are unreliable for defining clades.

225    Therefore, they should be analyzed with great caution.

226    Since Bayesian analysis allows to infer phylogenetic patterns from tree distributions, it

227    represents a more reliable tool to compare different evolutionary behaviors. Bayesian

228    analysis helps to obtain a tree topology that is closer to reality in the current conditions of

229    SARS-CoV-2 pandemic [26]. The phylogenetic analysis for nsp1, E, and Orf6 regions

230    confirmed the star-like topologies in accordance to a lower diversification of these regions

231    using the sequences available up to FN8 (Figure 2). Trees generated from nsp14 and Orf8

232    are at an intermediate point, where the formation of small clusters can be observed. In fact,

233    a mutation at position 28,144 (Orf8: L84S) has been proposed as a possible marker for viral

234    classification [27,28]. Finally, trees obtained from regions Orf3a, N, nsp3, and S showed the

235    best clade formation. Indeed, in the most variable regions nsp3 and S, it can be clearly seen

236    that sequences are separated into two large groups. Despite the aforementioned for the

237    nsp3 and S regions, even clusters with very high support values should be taken with

238    precaution and longer periods should be considered to obtain more accurate phylogeny

239    data. However, even when data are not the most accurate to study the spread or clade

240    formation [29, 30], they provide a good representation of the way the virus is evolving.

241    The analysis of amino acids frequencies allowed identifying different degree of region

242    conservation throughout the viral genome as a consequence of positive and negative

243    pressures. In particular, nsp3, S, Orf8, and N showed some substitutions in high

244    frequencies. This would indicate, as other authors previously report, the frequent circulation

245    of polymorphisms due to significant positive pressure [13,27,31]. Additionally, since S and N are

246    among candidates to be used in the formulation of vaccines and antibody treatment, it will be

247  important to monitor these substitutions in different geographic regions in order to improve

248  treatment and vaccination efficacy [32,33,34]. In particular, the appearance of the D614G variant

249  in the third week and its rapid increase until reaching a prevalence of 88% in the eighth week

250  could reflect an improvement in viral fitness, as several studies reported [35].

251  Contrarily, in regions nsp1, nsp14, E, and Orf6 no substitutions were selected and lasted

252  during the first 4 months of the pandemic. This would suggest that these are regions with

253  constraints to change due to the great negative selection pressure, as it has been recently

254  reported [13].

255  In the present study, the evolutionary rate for SARS-CoV-2 genes was estimated by

256  analyzing a large number of sequences, which were carefully curated and had a good

257  temporal and spatial structure. Additionally, the most phylogenetically informative regions of

258  the genome (nsp3 and S) were used for analysis, reinforcing the results confidence.

259  Previous studies on SARS-CoV-2 have reported similar data ranging from $1.79 \times 10^{-3}$ to

260  $6.58 \times 10^{-3}$ s/s/y for the complete genome [6,36]. However, in both articles, small datasets of

261  complete genomes were used (N=32 and 54, respectively). As studies were performed early

262  in the outbreak and due to datasets temporal structure, analysis could have led to less

263  precise estimates of the evolutionary rate [23]. On the other hand, another study from van

264  Dorp et al. (2020), analyzing 7,666 sequences has obtained different results with a

265  remarkably low evolutionary rate ($6 \times 10^{-4}$ nucleotide/genome/year) [15]. However, it is

266  important to consider that van Dorp et al. (2020) estimate the evolutionary rate using the

267  complete genome, including several highly conserved genomic regions, while in our work,

268  the estimation was performed with the most variable regions of the genome. Additionally,

269  tests randomizing the dates of nsp3 and S datasets were carried out; they showed that these

270  partial genomic regions have enough temporal signal. In this context, our results ($1.37 \times 10^{-3}$

271  s/s/y for NSp3 and $2.19 \times 10^{-3}$ s/s/y for S) are in close agreement with those published for

272 SARS-CoV genome, which have been estimated between 0.80 to 3.01 x $10^{-3}$ s/s/y [37-39](The

273 Chinese SARS Molecular Epidemiology Consortium, 2004, Vega et al. 2004, Zhao et al.

274 2004). Moreover, our values are in the same order magnitude as other RNA viruses [40].

275 Even though we should be cautious with these results interpretation, the date-randomization

276 analysis indicated a robust temporal signal.

277 In addition, the importance of separately studying the evolutionary rate in S region arises

278 from the fact that it represents the main target for antiviral agents and vaccines since it

279 includes the SARS-CoV-2 binding receptor domain (RBD), a crucial structure for the virus to

280 enter host cells and binding site for neutralizing antibodies [41].

281 Despite limitations of the evolutionary study of an emerging virus, where the selection

282 pressures are still low and therefore its variability is also low, this work has a great strength:

283 it lies on the extremely careful selection of a big sequence dataset to be analyze. First, it was

284 considered selected sequences to have a good temporal signal and spatial (geographic)

285 structure. Secondly, much attention was paid to the elimination of sequences with low

286 coverage and indeterminacies that could generate a noise for the phylogenetic analysis of a

287 virus that is beginning to evolve in a new host.

288 The appearance of a new virus means an adaptation challenge. The SARS-CoV-2 overcome

289 the spill stage and shows a significantly higher spread than SARS-CoV and MERS-CoV,

290 thus becoming itself the most important pandemic of the century. In this context, the results

291 obtained in this work about the variable diversity of nine crucial viral regions and the

292 determination of the evolutionary rate, are consequently decisive to understanding essential

293 feature of viral emergence. Nevertheless, monitoring SARS-CoV-2 population will be

294 required to determine the evolutionary course of new mutations as well as to understand the

295 way they affect viral fitness in human hosts.

296

297 **Competing interest**: On behalf of all authors, the corresponding author states that there is

298 no conflict of interest.

299

300 **Funding:** None

301

302 **Declaration of Author Contributions**

303 MJP: Data curation, acquisition of data, analysis and interpretation of data, drafting the

304 article, final approval of the version to be submitted.

305 LM: Data curation, acquisition of data, analysis and interpretation of data, revising the article

306 critically for important intellectual content, final approval of the version to be submitted.

307 APM: Data curation, Validation, revising the article critically for important intellectual content,

308 final approval of the version to be submitted.

309 DMF: Data curation, Validation, drafting the article, final approval of the version to be

310 submitted.

311 GG: Data curation, acquisition of data, analysis and interpretation of data, drafting the article,

312 final approval of the version to be submitted.

313 FAD: Conception and design of the study, acquisition of data, analysis and interpretation of

314 data, drafting the article, final approval of the version to be submitted.

315

316 **Acknowledgements**

317 MJP, LM, DMF, and FAD are members of the National Research Council (CONICET). We

318 would like to thank to the researchers who generated and shared the sequencing data from

319 GISAID (https://www.gisaid.org/) and Mrs. Silvina Heisecke from CEMIC-CONICET for

320 providing language assistance.

321

322 **REFERENCES**

323 [1] Su S, Wong G, Shi W, et al. Epidemiology, genetic recombination, and pathogenesis of

324 coronaviruses. *Trends in Microbiol*gy 2016; 24, 490-502.

325 https://doi.org/10.1016/j.tim.2016.03.003

326 [2] Cavanagh D. Coronavirus avian infectious bronchitis virus. *Veterinary Research* 2007*;*

327 38, 281-297. https://doi.org/10.1051/vetres:2006055

328 [3] Ismail MM, Tang AY & Saif YM. Pathogenicity of turkey coronavirus in turkeys and

329 chickens. *Avian Diseases* 2003; 47, 515-522. https://doi.org/10.1637/5917

330 [4] Zhu N, Zhang D, Wang W, et al. A Novel Coronavirus from Patients with Pneumonia in

331 China, 2019. *The New England Journal of Medicine* 2020; 382, 727-733.

332 https://doi.org/10.1056/NEJMoa2001017

333 [5] Coronaviridae Study Group of the International Committee on Taxonomy of Viruses. The

334 species *Severe acute respiratory syndrome-related coronavirus*: classifying 2019-nCoV and

335 naming it SARS-CoV-2. *Nature Microbiology 2020;* 5, 536-544.

336 https://doi.org/10.1038/s41564-020-0695-z

337 [6] Li X, Wang W, Zhao X, et al. Transmission dynamics and evolutionary history of 2019-

338 nCoV. *Journal of Medical Virology* 2020a; 92, 501-511. https://doi.org/10.1002/jmv.25701

339 [7] Li Q, Guan X, Wu P, et al. Early Transmission Dynamics in Wuhan, China, of Novel

340 Coronavirus-Infected Pneumonia. *The New England Journal of Medicine* 2020b; 382, 1199-

341 1207. https://doi.org/10.1056/NEJMoa2001316

342 [8] World Healt Organization, 2020. Coronavirus disease (COVID-19) Situation Report –

343 118. Retrieved from: https://www.who.int/docs/default-source/coronaviruse/situation-

344 reports/20200517-covid-19-sitrep-118.pdf?sfvrsn=21c0dafe_6 (15 August 2020, date last

345 accessed).

346 [9] Cui J, Li F & Shi ZL. Origin and evolution of pathogenic coronaviruses. *Nature Reviews*

347 *Microbiology* 2019; 17, 181-192. https://doi.org/10.1038/s41579-018-0118-9

348 [10] Luk HKH, Li X, Fung J, et al. Molecular epidemiology, evolution and phylogeny of SARS

349 coronavirus. *Infection Genetics and Evolution* 2019; 71, 21-30.

350 https://doi.org/10.1016/j.meegid.2019.03.001

351 [11] Zhou P, Yang XL, Wang XG, et al. A pneumonia outbreak associated with a new

352 coronavirus of probable bat origin. *Nature* 2020; 579, 270-273.

353 https://doi.org/10.1038/s41586-020-2012-7

354 [12] Benvenuto D, Giovanetti M, Salemi M, et al. The global spread of 2019-nCoV: a

355 molecular evolutionary analysis. *Pathogens and Global Health 2020;* 114, 64-67.

356 https://doi.org/10.1080/20477724.2020.1725339

357 [13] Cagliani R, Forni D, Clerici M, et al. Computational inference of selection underlying the

358 evolution of the novel coronavirus, SARS-CoV-2. *Journal of Virology* 2020;

359 https://doi.org/10.1128/JVI.00411-20

360 [14] Phan T. Genetic diversity and evolution of SARS-CoV-2. *Infection Genetics and*

361 *Evolution* 2020; 81, 104260.  https://doi.org/10.1016/j.meegid.2020.104260

362 [15] van Dorp L, Acman M, Richard D, et al. Emergence of genomic diversity and recurrent

363 mutations in SARS-CoV-2. *Infection Genetics and Evolution 2020;* 5, 104351.

364 https://doi.org/10.1016/j.meegid.2020.104351

365 [16] Strimmer K & von Haeseler A. Likelihood-mapping: A simple method to visualize

366 phylogenetic content of a sequence alignment. *Proceedings of the National Academy of*

367 *Sciences of the USA* 1997; 94, 6815-6819. https://doi.org/10.1073/pnas.94.13.6815

368 [17] Schmidt HA, Strimmer K, Vingron M, et al. TREE-PUZZLE: Maximum likelihood

369 phylogenetic analysis using quartets and parallel computing. *Bioinformatics* 2002; 18, 502-

370 504. https://doi.org/10.1093/bioinformatics/18.3.502

371   [18] Darriba D, Taboada GL, Doallo R, et al. jModelTest 2: more models, new heuristics and

372   parallelcomputing. *Nature Methods* 2012; 9, 772. https://doi.org/10.1038/nmeth.2109

373   [19] Ronquist F, Teslenko M, van der Mark P, et al. MrBayes 3.2: efficient Bayesian

374   phylogenetic inference and model choice across a large model space. *Systematic Biology*

375   2012; 61, 539-542. https://doi.org/10.1093/sysbio/sys029

376   [20] Rambaut A, Drummond AJ, Xie D, et al. Posterior summarization in Bayesian

377   phylogenetics using Tracer 1.7. *Systematic Biology* 2018; 67, 901-904.

378   https://doi.org/10.1093/sysbio/syy032

379   [21] Suchard MA, Lemey P, Baele G, et al. Bayesian phylogenetic and phylodynamic data

380   integration using BEAST 1.10. *Virus Evolution* 2018; 4, vey016.

381   https://doi.org/10.1093/ve/vey016

382   [22] Miller MA, Pfeiffer X, & Schwartz T. Creating the CIPRES Science Gateway for

383   inference of large phylogenetic trees. *Gateway Computing Environments Workshop* 2010; 1-

384   8. https://doi.org/10.1109/GCE.2010.5676129

385   [23] Duchene S, Featherstone L, Haritopoulou-Sinanidou M, et al. Temporal signal and the

386   phylodynamic threshold of SARS-CoV-2. *bioRxiv* 2020; [Preprint].

387   https://doi.org/10.1101/2020.05.04.077735

388   [24] Grassly NC & Fraser C. Mathematical models of infectious disease transmission. *Nat*

389   *Rev Microbiol 2008;* 6, 477-487. https://doi.org/10.1038/nrmicro1845

390   [25] Rieux A & Khatchikian CE. tipdatingbeast: an r package to assist the implementation of

391   phylogenetic tip-dating tests using beast. *Molecular Ecology Resources* 2017; 17, 608-613.

392   https://doi.org/10.1111/1755-0998.12603

393   [26] Drummond AJ, Ho SY, Phillips MJ, et al. Relaxed phylogenetics and dating with

394   confidence. *PLoS Biology* 2006; 4, e88. https://doi.org/10.1371/journal.pbio.0040088

395   [27] Tang X, Wu C, Li X, et al. On the origin and continuing evolution of SARS-CoV-2.

396   *National Science Review* 2020; 0, 1-12. https://doi.org/10.1093/nsr/nwaa036

397   [28] Yin C. Genotyping coronavirus SARS-CoV-2: methods and implications. *Genomics*

398   2020; 30318-30319. https://doi.org/10.1016/j.ygeno.2020.04.016

399   [29] Mavian C, Marini S, Prosperi M, et al. A snapshot of SARS-CoV-2 genome availability

400   up to 30th March, 2020 and its implications. JMIR Public Health Surveill 2020; 6, e19170

401   https://doi.org/10.2196/19170

402   [30] Sánchez-Pacheco SJ, Kong S, Pulido-Santacruz P, et al. Median-joining network

403   analysis of SARS-CoV-2 genomes is neither phylogenetic nor evolutionary. *Proceedings of*

404   *the National Academy of Sciences of the USA* 2020; 117, 9241–9243.

405   https://doi.org/10.1073/pnas.2007062117

406   [31] Issa E, Merhi G, Panossian B, et al. S.SARS-CoV-2 and ORF3a: Nonsynonymous

407   Mutations, Functional Domains, and Viral Pathogenesis. *mSystems* 2020 [Preprint].

408   https://doi.org/10.1128/mSystems.00266-20

409   [32] Ahmed SF, Quadeer AA & McKay MR. Preliminary Identification of Potential Vaccine

410   Targets for the COVID-19 Coronavirus (SARS-CoV-2) Based on SARS-CoV Immunological

411   Studies. *Viruses 2020;* 12, 254. https://doi.org/10.3390/v12030254

412   [33] Callaway E. The race for coronavirus vaccines: a graphical guide. *Nature* 2020*;* 580,

413   576-577. https://doi.org/10.1038/d41586-020-01221-y

414   [34] Koyama T, Weeraratne D, Snowdon JL, et al. Emergence of Drift Variants That May

415   Affect COVID-19 Vaccine Development and Antibody Treatment. *Pathogens* 2020; 9, 324.

416   https://doi.org/10.20944/preprints202004.0024.v1

417   [35] Li Q, Wu J, Nie J, et al. The Impact of Mutations in SARS-CoV-2 Spike on Viral

418   Infectivity and Antigenicity. *Cell* 2020; S0092-8674(20)30877-1. Advance online publication.

419   https://doi.org/10.1016/j.cell.2020.07.012

420  [36] Giovanetti M, Benvenuto D, Angeletti S, et al. The first two cases of 2019-nCoV in Italy:

421  Where they come from? *Journal of Medical Virology* 2020; 92, 518-521.

422  https://doi.org/10.1002/jmv.25699

423  [37] The Chinese SARS Molecular Epidemiology Consortium. Molecular Evolution of the

424  SARS Coronavirus During the Course of the SARS Epidemic in China. *Science* 2004; 303,

425  1666-1669. https://doi.org/10.1126/science.1092002

426  [38] Vega VB, Ruan Y, Liu J, et al. Mutational dynamics of the SARS coronavirus in cell

427  culture and human populations isolated in 2003. *BMC Infectious Diseases* 2004, 4, 32.

428  https://doi.org/10.1186/1471-2334-4-32

429  [39] Zhao Z, Li H, Wu X, et al. Moderate mutation rate in the SARS coronavirus genome and

430  its implications. *BMC Evolutionary Biology* 2004; 4, 21. https://doi.org/10.1186/1471-2148-4-

431  21

432  [40] Sanjuán R. From molecular genetics to phylodynamics: evolutionary relevance of

433  mutation rates across viruses. *PLoS Pathogens* 2012; 8, e1002685. https://doi.org/

434  10.1371/journal.ppat.1002685

435  [41] Ju B, Zhang Q, Ge J, et al. Human neutralizing antibodies elicited by SARS-CoV-2

436  infection. *Nature* 2020; 115–119. https://doi.org/10.1038/s41586-020-2380-z

437

438

439

440

441

442

443

444

445    **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** |

446    FN: Fortnight; Q1=quartile 1, Q3=quartile 3. The total number of sequences is variable depending on
447    the analyzed region (nsp1, 1608; nsp3, 1511; nsp14, 1550; S, 1488; Orf3a, 1600; E, 1615; Orf6,
448    1616; Orf8, 1612; and N, 1610)
449

450 **Table 2.** Amino acids selected by region and fortnight. The number indicates the amino

451 acids 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 |
| | L84S | 0 | 42.0 | 37.0 | 21.0 | 21.0 | 18.0 | 7.0 | 6.0 |
| 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 |

452 Only regions where amino acid change was selected and remained until the last analyzed fortnight
453 are shown. FN: Fortnight; aa: amino acid

454

455

456 **Table 3.** Number of variable positions, number of mutations, and number of sequences with

457 mutation by region

| Region | Nº of variable aa positions (%) | Nº of aa substitutions | Nº 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 | - |

458 aa: amino acid

459

460

461

462

463

464

465

466

467

468

469

470

471 **FIGURE LEGENDS**

472

473 **Figure 1**

474 Phylogenetic signal for SARS-CoV-2 datasets. Presence of phylogenetic signal was

475 evaluated by likelihood mapping, unresolved quartets (center) and partly resolved quartets

476 (edges) for genomes available on April 17 for the nine analyzed regions: nsp1 (29

477 sequences), nsp3 (225 sequences), nsp14 (65 sequences), S (183 sequences), Orf3a (74

478 sequences), E (11 sequences), Orf6 (12 sequences), Orf8 (23 sequences), and N (113

479 sequences). Presence of strong phylogenetic signal (<40% unresolved quartets) was not

480 reached for any region.

481

482 **Figure 2**

483 Bayesian trees of 29 sequences of nsp1 (540nt), 225 sequences of nsp3 (5835nt), 65

484 sequences of nsp14 (1581nt), 183 sequences of S (3822nt), 74 sequences of Orf3a (828nt),

485 11 sequences of E (228nt), 12 sequences of Orf6 (186nt), 23 sequences of Orf8 (366nt),

486 and113 sequences of N (1260nt). Scale bar represents substitutions per site.

487

488 **Figure 3**

489 Comparison of the evolutionary rates estimated using BEAST for the original dataset and the

490 date-randomized datasets (312 sequences). This analysis was performed for regions nsp3

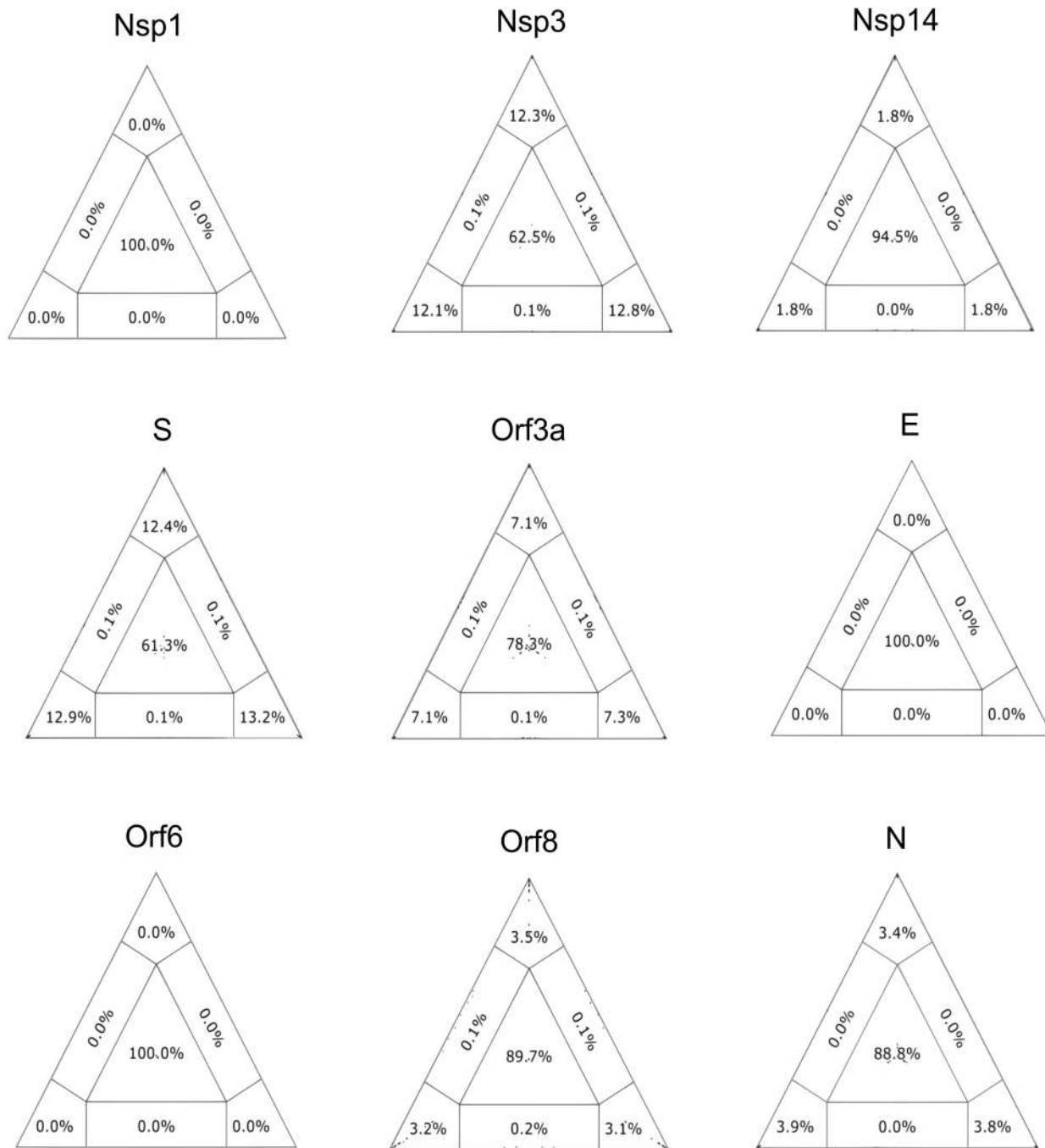491 (5835nt) and S (3822nt). s.s.y = substitutions/site/year.
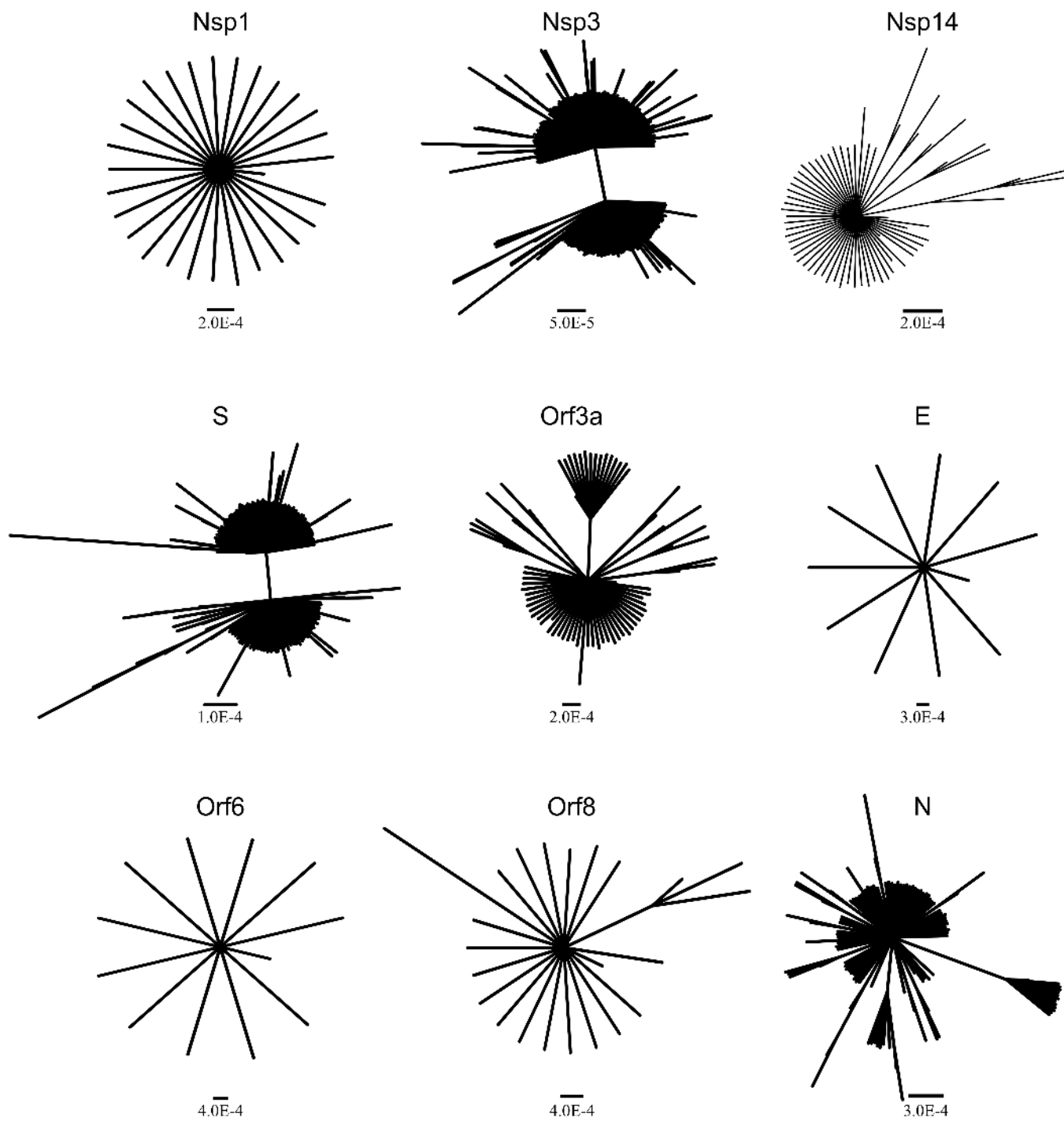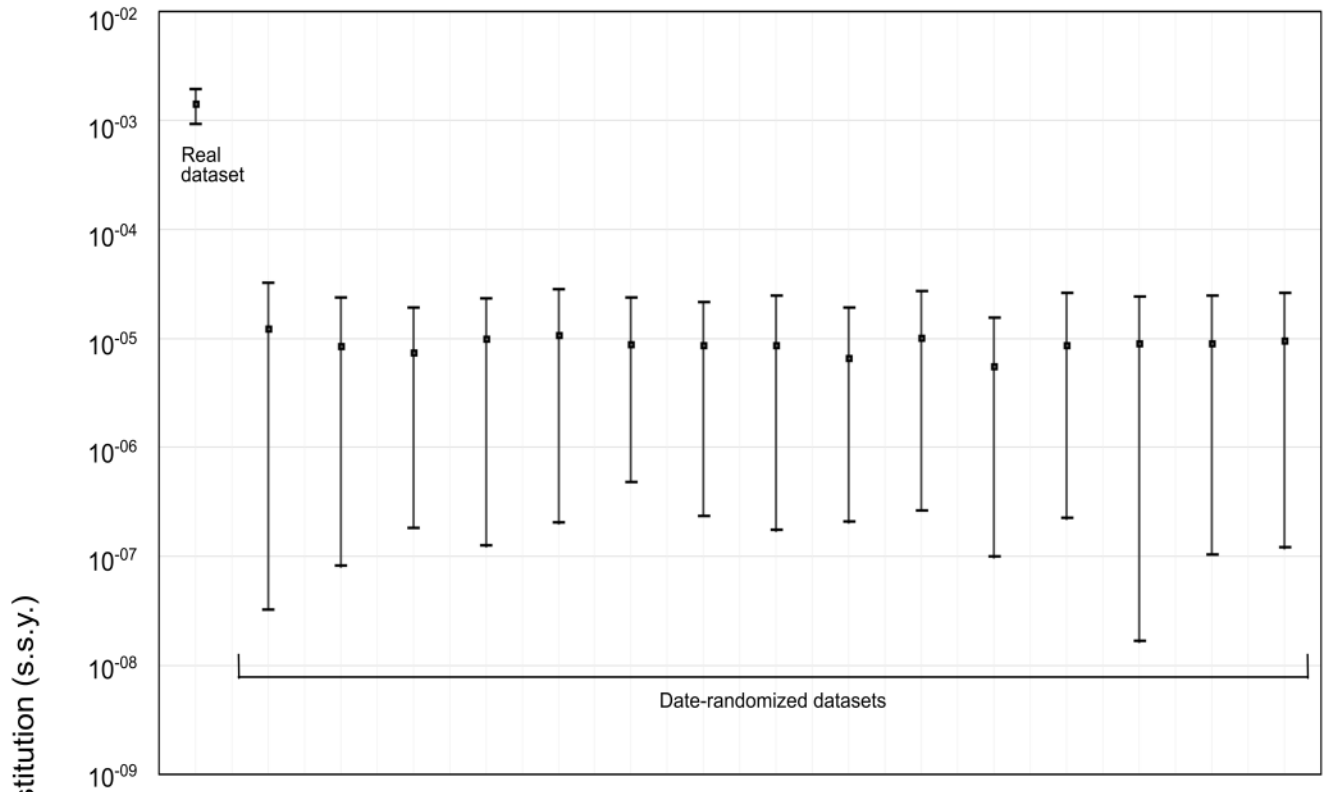
492

493

Figure 1

Nsp1

2.0E-4

Nsp3

5.0E-5

Nsp14

2.0E-4

S

1.0E-4

Orf3a

2.0E-4

E

3.0E-4

Orf6

4.0E-4

Orf8

4.0E-4

N

3.0E-4

Figure 2

Figure 3