

The first wave of the Spanish COVID-19 epidemic was associated with early introductions and fast spread of a dominating genetic variant

Mariana G. López^{1,#}, Álvaro Chiner-Oms^{1,#}, Darío García de Viedma^{2,3,4}, Paula Ruiz-Rodriguez⁵, María Alma Bracho^{6,7}, Irving Cancino-Muñoz¹, Giuseppe D'Auria^{7,8}, Griselda de Marco⁸, Neris García-González⁶, Galo Adrian Goig⁹, Inmaculada Gómez-Navarro¹, Santiago Jiménez-Serrano¹, Lúcia Martínez-Priego⁸, Paula Ruiz-Hueso⁸, Lidia Ruiz-Roldán⁶, Manuela Torres-Puente¹, Juan Alberola^{10,11,12}, Eliseo Albert¹³, Maitane Aranzamendi Zaldumbide^{14,15}, María Pilar Bea-Escudero¹⁶, Jose Antonio Boga^{17,18}, Antoni E. Bordoy¹⁹, Andrés Canut-Blasco²⁰, Ana Carvajal²¹, Gustavo Cilla Eguiluz²², Maria Luz Cordón Rodríguez²⁰, José J. Costa-Alcalde²³, María de Toro¹⁶, Inmaculada de Toro Peinado²⁴, Jose Luis del Pozo²⁵, Sebastián Duchêne²⁶, Jovita Fernández-Pinero²⁷, Begoña Fuster Escrivá^{28,29}, Concepción Gimeno Cardona²⁸, Verónica González Galán^{30,31}, Nieves Gonzalo Jiménez³², Silvia Hernáez Crespo²⁰, Marta Herranz^{2,3,4}, José Antonio Lepe^{30,31}, José Luis López-Hontangas³³, María Ángeles Marcos³⁴, Vicente Martín^{7,35}, Elisa Martró^{7,19}, Ana Milagro Beamonte³⁶, Milagrosa Montes Ros²², Rosario Moreno-Muñoz³⁷, David Navarro^{13,29}, José María Navarro-Mari^{38,39}, Anna Not¹⁹, Antonio Oliver^{40,41}, Begoña Palop-Borrás²⁴, Mónica Parra Grande⁴², Irene Pedrosa-Corral^{38,39}, Maria Carmen Perez Gonzalez⁴³, Laura Pérez-Lago^{2,3}, Luis Piñeiro Vázquez²², Nuria Rabella^{44,45,46}, Jordi Reina⁴⁰, Antonio Rezusta^{36,47,48}, Lorena Robles Fonseca⁴², Ángel Rodríguez-Villodres^{30,31}, Sara Sanbonmatsu-Gámez^{38,39}, Jon Sicilia^{2,3}, María Dolores Tirado Balaguer³⁷, Ignacio Torres¹³, Alexander Tristancho^{36,47}, José María Marimón²², Mireia Coscolla^{5,*}, Fernando González-Candelas^{6,7,*}, Iñaki Comas^{1,7,*} on behalf of the SeqCOVID-SPAIN consortium.

SUPPLEMENTARY NOTES

Large metropolitan areas as sources of viral genetic diversity

We observed a heterogeneous distribution of SARS-CoV-2 genetic diversity across Spain, both at regional and local levels (Figure 1b). Some regions concentrate different viral lineages while others have low genetic diversity and are dominated by one or few lineages. We found a weak correlation between the mean viral genetic diversity and the population density of each

35 municipality ($\rho=0.32$, $p\text{-value}=0.01$), suggesting that more densely populated areas exhibit higher
36 diversity of SARS-CoV-2 strains compared to low-density regions.

37 We hypothesize that large cities acted as sources of viral genetic diversity during the pandemic,
38 with their “economic-influence” areas (e.g. urban outdoor areas) acting as sinks. To test this, we
39 used the autonomous region of Comunidad Valenciana as a case study, since this was the region
40 for which we have the largest sampling. We calculated the genetic diversity and the mean pairwise
41 distances between samples belonging to different municipalities and correlated these values
42 against their geographic distance to the regional capital city, Valencia. Results (Figure S10, Figure
43 S11) suggest that the city of Valencia was the center of the genetic diversity within the region; the
44 genetic diversity decreases with geographic distance to Valencia up to 70-80km away from the city
45 center. This result agrees with a “model of isolation” by distance, where only part of the genetic
46 diversity of the original population (Valencia) spreads to the surrounding areas; it is also
47 compatible with patterns of population mobility within large cities before the installment of a
48 national lockdown.

49

50 Spread of SECs across the country

51 We observed a correlation between the time of introduction and SEC size (Figure S6). In the
52 beginning, before the adoption of non-pharmaceutical interventions, strains were transmitted
53 without control and generated larger clusters. After restrictions were implemented, transmission
54 was controlled and smaller clusters were observed. Regardless of their introduction time, some
55 strains could spread locally while others dispersed throughout the rest of Spain, generating
56 different dispersion patterns. To study this, we first analyzed the geographic distribution of
57 samples within each SEC (Figure 2d). We noticed that the largest and oldest SECs (SEC7 and
58 SEC8) were the most widely distributed and were present in 10 or more Spanish regions.
59 Conversely, the smallest and youngest SECs (SEC1 and SEC2) spread only locally, in Andalucía
60 and País Vasco respectively, in agreement with its MRCA ancestor occurring after the lockdown.
61 The remaining SECs displayed varying but generally narrower dispersal patterns compared to
62 SEC7 and SEC8, and are present in 2 up to 6 Spanish regions. We also analyzed the distribution
63 of geographic distances (kilometers) between samples within the same SEC. An ANOVA analysis
64 (adjusted $p\text{-value} \ll 0.01$) divides SECs into three different categories: SECs 1,2,4,6 and 9
65 showed very limited geographical dispersion; SEC3 and 5 showed medium values and SEC7 and
66 8 exhibited the widest geographical distribution. The smallest SECs include samples within a
67 range of about 0-58 km (Figure S7), whereas the mean distance between samples in SEC7 is
68 323 km (interquartile range 225-447 km) and 371 km (interquartile range 76-634 km) in SEC8.
69 This supports the hypothesis that SECs 7 and 8 spread rapidly throughout the country after their
70 first introduction. Interestingly, despite their wider geographic range compared to the rest of SECs,
71 they do not exhibit a larger accumulation of mutations (Figure S7).

72

73 Impact of lockdown

74 In order to evaluate the effectiveness of the restriction movements on the spread of the pandemic
75 in Spain, we performed a Bayesian birth-death skyline (BDSKY) analysis in order to estimate the

76 timing and magnitude of changes in the effective reproductive number (R_e), as a measure of the
77 average secondary cases generated by an infected person, for the most successful SECs. For
78 SEC7 the results suggested a strong evidence for a change in R_e around 20th March with a
79 decrease below one, after the restrictions implemented by the Spanish Government on 14th
80 March(Figure S9).

81 In the case of the SEC8, the decrease of the R_e was estimated after 9th March, a bit earlier than
82 the lockdown implementation, this is probably due to the bias produced by the multiple
83 introductions that cannot be differentiated from transmission events by the birth-death skyline
84 model as can be observed in the multiple and basal tree bifurcations (Figure 3c).

85 Doubling time, as measure of the amount of time in which the incidence doubles ¹, was also
86 calculated for both SECs as $365 \text{ days/year} \times \ln(2) / (R_e \times \delta) - \delta$, where δ is known as the become
87 uninfected rate and is the inverse of the duration of infection, which we assume to be 10 days,
88 or 10/365 years (such that $\delta=36.5 \text{ year}^{-1}$). Before the corresponding R_e changing date, the
89 doubling time was estimated at 6.3 days (95% HPD: 4.3-10.2 days) for SEC7 and 3.3 days (95%
90 HPD: 2.7 - 4.1 days) for SEC8. After changing date R_e value is below 1, the number of cases is
91 decreasing, such that the absolute value of the doubling time is in fact the halving time of the
92 number of infections, i.e the time required for the number of cases to halve. The post-peak halving
93 times are approximately 9.5 and 9.0 days for SEC7 and SEC8 respectively, reinforcing the results
94 of a pandemic reduction after the confinement.

95

96 SUPPLEMENTARY MATERIALS AND METHODS

97 SeqCOVID sampling and sequencing

98 RNA samples were received from different hospitals, and confirmed as SARS-Cov-2-positive by
99 RT-PCR by Microbiological Services. Samples were the remaining RNA extracts from naso- and
100 oropharyngeal clinical specimens employed for diagnosis. The use of such samples have been
101 approved by the ethics committee Comité Ético de Investigación de Salud Pública y Centro
102 Superior de Investigación en Salud Pública (CEI DGSP-CSISP) N° 20200414/05. RNA was retro-
103 transcribed into cDNA and SARS-CoV-2 complete genome amplification was conducted in two
104 multiplex PCR, accordingly to openly available protocol developed by the ARTIC network ² using
105 the V3 multiplex primers scheme ³. Two resulting amplicon pools were combined and used for
106 library preparation. Genomic libraries were constructed with the Nextera DNA Flex Sample
107 Preparation kit (Illumina Inc., San Diego, CA) according to the manufacturer's protocol, with 5
108 cycles for indexing PCR. Whole genome sequencing was carried out in the MiSeq platform (2×200
109 cycles paired-end run; Illumina).

110 The sequences obtained went through a bioinformatic pipeline based on IVAR ⁴, which is open
111 source and can be accessed at <https://gitlab.com/fisabio-ngs/sars-cov2-mapping>. In short, the
112 pipeline goes through the following steps: 1) Removal of the human reads with Kraken ⁵; 2)
113 filtering of the fastq files using fastp v 0.20.1 ⁶ (arguments: --cut tail, --cut-window-size, --cut-
114 mean-quality, -max_len1, -max_len2); 3) mapping and variant calling using IVAR v 1.2; 4) quality
115 control assessment with MultiQC ⁷.

116

117 Global alignment and phylogenetic reconstruction

118 To build the global alignment, we downloaded and concatenated all non-spanish sequences
119 present in GISAID ⁸ on 21st June that passed strict filtering criteria: i) sequences should have
120 more than 29,000 bp length, ii) verified insertions/deletions, iii) less than 1% of Ns and less than
121 0.05% of unique amino acid mutations (compared with other sequences in GISAID).

122 Later, we added all Spanish sequences deposited in GISAID up to July 29th. The final alignment
123 constructed included 32,914 sequences. The accession numbers of the sequences used in this
124 study can be found in Table S1.

125 Sequences were aligned against the SARS-CoV-2 reference genome ⁹ using MAFFT ¹⁰. Specific
126 positions that have been reported to be problematic for phylogenetic reconstruction ¹¹ were
127 masked, following the procedure described by Rob Lanfear ¹², using the mask_alignment.sh
128 script.

129 Finally, a maximum-likelihood (ML) phylogeny was reconstructed using IQTREE ¹³ with the GTR
130 model and based on the complete masked genome alignment. This phylogeny was rooted to the
131 SARS-CoV-2 sequence obtained in Wuhan on 24/12/2019 (GISAID ID: EPI_ISL_402123).
132

133 Identification of introductions and transmission clusters

134 We identified transmission groups between Spanish sequences by inspecting the global
135 phylogeny (32,914 leaves) and searching for Spanish sequences (or groups of) that were
136 embedded within sequences with other geographic origins. Given the general low diversity among
137 sequences, most phylogenetic nodes ended up being polytomic in the maximum-likelihood tree.
138 Because of this, we defined three different transmission scenarios: i) strains that represent
139 introductions in Spain but differ from those from other countries and form well defined
140 transmission groups ('candidate transmission clusters'); ii) strains introduced into Spain that are
141 equal to other Spanish sequences and that are also equal to sequences from other countries
142 ('zero distance clusters'); and iii) Spanish sequences found within groups of sequences from other
143 countries and which are not phylogenetically near any other Spanish sequences ('unique'). The
144 'candidate transmission clusters' were identified as monophyletic groups of sequences composed
145 exclusively by Spanish sequences in the phylogeny. The 'zero distance clusters' were identified
146 as Spanish sequences that share a common ancestor and that are at 0 SNP distance from each
147 other. Finally, The 'unique' sequences were identified as those sequences which do not share
148 their most recent ancestor with any other Spanish sequence.

149 Next, we inferred how many of these transmission groups have a potential contagion date for their
150 first case that predates the start of mobility restrictions, on 14th March, by subtracting 14 days to
151 the diagnosis date.

152 Finally, we wanted to investigate the international origin of these introductions. For each of the
153 identified groups or 'unique' sequences with an inferred contagion date before 14th March, we
154 looked for the closest non-Spanish sequence in the phylogeny with a diagnosis date predating
155 the first case of the transmission group. As the current consensus is that the pandemic began in
156 Asia and later it moved to Europe, we considered only those sequences with an Asian or
157 European origin as potential sources of introductions.
158

159 SEC alignment and phylogeny

160 Using the global phylogeny, we identified nodes which had at least 20 leaves and in which at least
161 50% of these correspond to Spanish sequences. Next, for each of these nodes or clades we
162 reconstructed an alignment of the complete masked genomes including:

- 163 a. The sequences that belong to the identified clade
- 164 b. 11 basal sequences from Wuhan acting as an 'anchor' for the phylogeny (Table
165 S1).
- 166 c. A subset of 51 representative sequences, each one from a different pangolin
167 lineage, selected to maximize the global SARS-CoV-2 genetic diversity (Table S1,
168 downloaded from GISAID on 2020-07-20).

169 For each of these alignments we inferred a ML phylogeny, using IQTREE ¹³, with the model GTR,
170 1,000 fast-bootstrap replicates and rooted to the Wuhan sequence (EPI_ISL_402123). Then, in
171 the resulting phylogeny we identified less inclusive nodes embedded within the above identified
172 clades and had a bootstrap support value > 80. These clades were named as potential Spanish
173 epidemic clades (SEC). The iTOL tool ¹⁴ was used for phylogenetic visualization.

174

175 SEC8 detailed analysis

176 To get more detail on SEC8 phylogenetic structure, and to evaluate if mobility restrictions were
177 effective to hinder SEC8 transmission, we enriched the original SEC8 phylogenetic tree with all
178 the isolates of this clade sampled from February to October by the SeqCOVID consortium (959
179 sequences in total). Later, epidemiological information was included and plotted in the tree using
180 the iTOL tool ¹⁴.

181

182 SEC8 potential superspreading events were defined as groups of more than 10 sequences,
183 having at least 1 SNP in common and having a within sequence median distance from 1 to 3
184 SNPs.

185

186

187 Population genetics and differentiation geography

188 Geographic distance between sequences were computed using the GPS coordinates of the
189 patient residence city and applying the Vicenty (ellipsoid) method. Genetic diversity was
190 calculated with two different methods: i) genetic distance between each pair of samples in number
191 of substitutions (SNPs), and ii) number of base substitutions per site averaged over all sequence
192 pairs in a group of sequences. Both values have been estimated using the MEGA software ¹⁵,
193 skipping one position when a gap is found in the two compared sequences.

194 Demographic data for all Spanish regions and municipalities were downloaded from INE
195 (<https://www.ine.es/>), and had been updated on 1st January, 2020.

196

197 The genetic diversity heatmap of the Comunidad Valenciana autonomous was generated with
198 QGIS v.3.14.16-Pi ¹⁶, using the inverse distance weighting (IDW) algorithm to interpolate the
199 mean genetic diversity of each municipality for which we had at least two sequences.

200
201 To compare the genetic and geographic distance distribution between the different SECs, we
202 used a one-way ANOVA test, followed by multiple pairwise-comparisons of the between-groups
203 mean with a Tukey HSD test.

204

205 Dating analyses

206 To estimate the most recent common ancestor (MRCA) of each of the nine SECs defined above,
207 a multi-sequence alignment was performed including the 11 samples belonging to basal
208 phylogenetic clades and the 51 representative sequences from different lineages (Table S1).
209 Before phylogenetic dating, root-to-tip regression of genetic divergence against sampling dates
210 was performed to investigate the molecular clock signal of SECs using TempEst v 1.5.3 ¹⁷. We
211 implemented a coalescent Bayesian exponential growth model available in Beast 2.6 ¹⁸ with the
212 HKY+ Γ model of nucleotide substitution. Tree priors were defined as follows: for effective
213 population size we used a lognormal distribution (M=1, S=2) and for growth rate a Laplace
214 distribution (M=0, S=100). The uncorrelated lognormal relaxed clock was selected as the best
215 fitting clock model using Bayes Factor comparisons of strict and relaxed clocks based on path
216 sampling/stepping stone analysis ¹⁹. Clock priors were defined as: uclid.mean: lognormal
217 distribution with mean in real space = 1.4×10^{-3} subs/site/year and uclid.stdev = 5×10^{-2} .
218 Parameters were estimated using Markov Chain Monte Carlo (MCMC) Bayesian inference, with
219 5×10^7 steps-long chains with exception of SEC7 and SEC8, for which longer chains were run (1
220 $\times 10^8$), in all cases a total of 10^5 steps were sampled in the log files. For all analysis, three
221 independent runs starting from different seeds were conducted in order to ensure convergence,
222 then combined with LogCombiner v 2.6.3 after removing the initial 10% of the MCMC as burn-in.
223 Adequate mixing of parameters and convergence among runs were assessed using Tracer v
224 1.7.1 ²⁰ by verifying that each parameter reached an effective sampling size (ESS) above 200
225 and that traces showed stationarity and good mixing. The final posterior distribution contained a
226 total of 9000 trees, annotated with Treeannotator v 2.6.3 and visualized in FigTree v 1.4.3 ²¹

227

228 Phylodynamics analysis to estimate R_e

229 To estimate discrete changes in R_e through for the two largest epidemic clades SEC7 and SEC8.
230 We used a Bayesian birth-death skyline model (BDSKY) with serial sampling ²² implemented in
231 BEAST v 2.6 ¹⁸. BDSKY uses an episodic, piecewise birth-death model in which the parameter is
232 allowed to change at discrete points in time, with the magnitude and timing of changes estimated
233 from the data. In our analysis, we set two intervals wherein R_e is constant and estimated the date
234 with most evidence for a change in R_e . To this end, we set a uniform prior distribution. R_e was
235 estimated before and after the changing time. Same parameters as above were used but fixing
236 the clock rate and the become uninfected rate ($\delta = 36.5 \text{ years}^{-1}$) in accordance with consistent
237 global estimates of an infectious period of 10 days ²³. In order to avoid bias in the model

238 parameters due to constant sampling proportion assumed by BDSKY models, this parameter was
239 set to zero before the first sample date using TreeSlicer
240 (<https://github.com/laduplessis/skylinetools/wiki>). For this analysis a 1×10^7 and 4×10^7 steps-
241 long chains were used for SEC7 and SEC8 respectively. Results were inspected with Tracer (v
242 1.7.1)²⁰ by verifying that every parameter had effective sampling sizes above 200 and well mixing
243 was obtained. Doubling time was calculated from the parameters estimated by BDSKY model in
244 which growth rate(r) = $(R_e * \delta) - \delta$ and doubling time = $\ln(2) / r$.

245 Statistical analysis

246 All statistical analyses were carried out using the R statistical language²⁴. Packages ape²⁵, treeio
247^{25,26}, doParallel²⁷ and foreach²⁸ were used for phylogenetic manipulation and analysis. We
248 additionally used packages geosphere²⁹, lwgeom³⁰, sp³¹, sf³² and rgeos³³ to calculate the
249 geographic distances between samples and the geographical representation in the data. The
250 ggplot2 R package³⁴ was extensively used for analysis and data plotting.
251

252 Data availability

253 The analysis pipeline used to map and analyze the sequences is available at
254 <https://gitlab.com/fisabio-ngs/sars-cov2-mapping>. All the genomic sequences used in the
255 analyses are available in the GISAID database, and the accession numbers can be found in Table
256 S1.

257 SUPPLEMENTARY REFERENCES

258

- 259 1. Lurie, M. N., Silva, J., Yorlets, R. R., Tao, J. & Chan, P. A. Coronavirus Disease 2019
260 Epidemic Doubling Time in the United States Before and During Stay-at-Home Restrictions.
261 *J. Infect. Dis.* **222**, 1601–1606 (2020).
- 262 2. Quick, J. *nCoV-2019 sequencing protocol*. [https://www.protocols.io/view/ncov-2019-
263 sequencing-protocol-bbmui6w.pdf](https://www.protocols.io/view/ncov-2019-sequencing-protocol-bbmui6w.pdf) (2020) doi:10.17504/protocols.io.bbmui6w.
- 264 3. artic-network. artic-network/artic-ncov2019. <https://github.com/artic-network/artic-ncov2019>.
- 265 4. Grubaugh, N. D. *et al.* An amplicon-based sequencing framework for accurately measuring
266 intrahost virus diversity using PrimalSeq and iVar. *Genome Biol.* **20**, 1–19 (2019).
- 267 5. Wood, D. E. & Salzberg, S. L. Kraken: ultrafast metagenomic sequence classification using
268 exact alignments. *Genome Biol.* **15**, 1–12 (2014).

- 269 6. Chen, S., Zhou, Y., Chen, Y. & Gu, J. fastp: an ultra-fast all-in-one FASTQ preprocessor.
270 *Bioinformatics* **34**, i884–i890 (2018).
- 271 7. Ewels, P., Magnusson, M., Lundin, S. & Källér, M. MultiQC: summarize analysis results for
272 multiple tools and samples in a single report. *Bioinformatics* **32**, 3047–3048 (2016).
- 273 8. Shu, Y. & McCauley, J. GISAID: Global initiative on sharing all influenza data - from vision
274 to reality. *Euro Surveill.* **22**, (2017).
- 275 9. Wu, F. *et al.* A new coronavirus associated with human respiratory disease in China.
276 *Nature* **579**, (2020).
- 277 10. Katoh, K., Misawa, K., Kuma, K.-I. & Miyata, T. MAFFT: a novel method for rapid multiple
278 sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* **30**, 3059–3066
279 (2002).
- 280 11. De Maio, N. *et al.* Issues with SARS-CoV-2 sequencing data. (2020).
- 281 12. Lanfear, R. *A global phylogeny of SARS-CoV-2 sequences from GISAID.* (Zenodo, 2020).
282 doi:10.5281/ZENODO.3958883.
- 283 13. Nguyen, L.-T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: a fast and effective
284 stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* **32**,
285 268–274 (2015).
- 286 14. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v4: recent updates and new
287 developments. *Nucleic Acids Res.* **47**, W256–W259 (2019).
- 288 15. Kumar, S., Stecher, G., Peterson, D. & Tamura, K. MEGA-CC: computing core of molecular
289 evolutionary genetics analysis program for automated and iterative data analysis.
290 *Bioinformatics* **28**, 2685–2686 (2012).
- 291 16. QGIS Development Team. *QGIS Geographic Information System.* (2020).
- 292 17. Rambaut, A., Lam, T. T., Max Carvalho, L. & Pybus, O. G. Exploring the temporal structure
293 of heterochronous sequences using TempEst (formerly Path-O-Gen). *Virus Evol* **2**, vew007
294 (2016).

- 295 18. Bouckaert, R. *et al.* BEAST 2.5: An advanced software platform for Bayesian evolutionary
296 analysis. *PLoS Comput. Biol.* **15**, e1006650 (2019).
- 297 19. Grummer, J. A., Bryson, R. W. & Reeder, T. W. Species delimitation using Bayes factors:
298 simulations and application to the *Sceloporus scalaris* species group (Squamata:
299 Phrynosomatidae). *Syst. Biol.* **63**, (2014).
- 300 20. Rambaut, A., Drummond, A. J., Xie, D., Baele, G. & Suchard, M. A. Posterior
301 Summarization in Bayesian Phylogenetics Using Tracer 1.7. *Systematic Biology* **67**, 901–
302 904 (2018).
- 303 21. Rambaut, A. *FigTree*. (2016).
- 304 22. Stadler, T., Kühnert, D., Bonhoeffer, S. & Drummond, A. J. Birth–death skyline plot reveals
305 temporal changes of epidemic spread in HIV and hepatitis C virus (HCV). *Proc. Natl. Acad.*
306 *Sci. U. S. A.* **110**, 228–233 (2013).
- 307 23. He, X. *et al.* Temporal dynamics in viral shedding and transmissibility of COVID-19. *Nat.*
308 *Med.* **26**, 672–675 (2020).
- 309 24. R Core Team. *R: A language and environment for statistical computing*. (2017).
- 310 25. Paradis, E. & Schliep, K. ape 5.0: an environment for modern phylogenetics and
311 evolutionary analyses in R. *Bioinformatics* **35**, 526–528 (2019).
- 312 26. Wang, L.-G. *et al.* Treeio: An R Package for Phylogenetic Tree Input and Output with Richly
313 Annotated and Associated Data. *Molecular Biology and Evolution* **37**, 599–603 (2020).
- 314 27. Microsoft Corporation and Steve Weston. *doParallel: Foreach Parallel Adaptor for the*
315 *'parallel' Package*. (Comprehensive R Archive Network (CRAN), 2019).
- 316 28. Microsoft Corporation and Steve Weston. *foreach: Provides Foreach Looping Construct*.
317 (Comprehensive R Archive Network (CRAN), 2020).
- 318 29. Hijmans, R. J., Williams, E. & Vennes, C. *Spherical Trigonometry*. (Comprehensive R
319 Archive Network (CRAN), 2019).
- 320 30. Pebesma, E., Rundel, C. & Teucher, A. *lwgeom: Bindings to Selected 'liblwgeom'*

321 *Functions for Simple Features*. (2020).

322 31. Bivand, R. S., Pebesma, E. & Gómez-Rubio, V. *Applied Spatial Data Analysis with R*.
323 (Springer Science & Business Media, 2013).

324 32. Pebesma, E. Simple Features for R: Standardized Support for Spatial Vector Data. *The R*
325 *Journal* **10**, 439 (2018).

326 33. Bivand, R. *et al.* *rgeos: Interface to Geometry Engine - Open Source ('GEOS')*.
327 (Comprehensive R Archive Network (CRAN), 2020).

328 34. Wilkinson, L. ggplot2: Elegant Graphics for Data Analysis by WICKHAM, H. *Biometrics* vol.
329 67 678–679 (2011).

330
331