



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology

Citation for published version:

Rambaut, A, Holmes, EC, O'Toole, Á, Hill, V, McCrone, JT, Ruis, C, du Plessis, L & Pybus, OG 2020, 'A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology', *Nature Microbiology*, vol. 5, pp. 1403-1407. <https://doi.org/10.1038/s41564-020-0770-5>

Digital Object Identifier (DOI):

[10.1038/s41564-020-0770-5](https://doi.org/10.1038/s41564-020-0770-5)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

Nature Microbiology

Publisher Rights Statement:

This is the accepted version of the following article: Rambaut, A., Holmes, E.C., O'Toole, Á. et al. A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology. *Nat Microbiol* (2020). <https://doi.org/10.1038/s41564-020-0770-5>, which has been published in final form at <https://www.nature.com/articles/s41564-020-0770-5>.

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



1 **A dynamic nomenclature proposal for SARS-CoV-2 lineages to**
2 **assist genomic epidemiology**

3

4

5 Andrew Rambaut^{1*}, Edward C. Holmes^{2*}, Áine O'Toole¹, Verity Hill¹, John T. McCrone¹,

6 Christopher Ruis⁴, Louis du Plessis³, Oliver G. Pybus^{3*}

7

8

9 ¹Institute of Evolutionary Biology, University of Edinburgh, Edinburgh, UK.

10 ²Marie Bashir Institute for Infectious Diseases and Biosecurity, School of Life and
11 Environmental Sciences and School of Medical Sciences, University of Sydney, Sydney, NSW,
12 Australia.

13 ³Department of Zoology, University of Oxford, Oxford, UK

14 ⁴Department of Medicine, University of Cambridge, UK.

15

16 *Correspondence should be addressed to:

17 a.rambaut@ed.ac.uk, edward.holmes@sydney.edu.au, oliver.pybus@zoo.ox.ac.uk

18

19

20 **The ongoing pandemic spread of a novel human coronavirus, SARS-COV-2, associated**
21 **with severe pneumonia disease (COVID-19), has resulted in the generation of tens of**
22 **thousands of virus genome sequences. The rate of genome generation is unprecedented, yet**
23 **there is currently no coherent nor accepted scheme for naming the expanding phylogenetic**
24 **diversity of SARS-CoV-2. We present a rational and dynamic virus nomenclature that uses**
25 **a phylogenetic framework to identify those lineages that contribute most to active spread.**
26 **Our system is made tractable by constraining the number and depth of hierarchical lineage**
27 **labels and by flagging and de-labelling virus lineages that become unobserved and hence**
28 **are likely inactive. By focusing on active virus lineages and those spreading to new**
29 **locations this nomenclature will assist in tracking and understanding the patterns and**
30 **determinants of the global spread of SARS-CoV-2.**

31

32

33 There are currently more than 35,000 publicly available complete or near-complete genome
34 sequences of SARS-CoV-2 (as of 1st June 2020) and the number continues to grow. This
35 remarkable achievement has been made possible by the rapid genome sequencing and online
36 sharing of SARS-CoV-2 genomes by public health and research teams worldwide. These
37 genomes have the potential to provide invaluable insights into the ongoing evolution and
38 epidemiology of the virus during the pandemic, and will likely play an important role in
39 surveillance and its eventual mitigation and control. Despite such a wealth of data, there is
40 currently no coherent system for naming and discussing the growing number of phylogenetic
41 lineages that comprise the population diversity of this virus, with conflicting *ad hoc* and informal
42 systems of virus nomenclature in circulation. A nomenclature system for the genetic diversity of
43 SARS-CoV-2 (a clade within the species *Severe acute respiratory syndrome-related virus*, sub-
44 genus *Sarbecovirus*, genus *Betacoronavirus*, family *Coronaviridae*¹) is urgently required before
45 scientific literature and communication become further confused.

46

47 There is no universal approach to classifying virus genetic diversity below the level of a virus
48 species², and this is not covered by the International Committee on Taxonomy of Viruses
49 (ICTV). Typically, genetic diversity is categorised into distinct ‘clades’, each of which
50 corresponds to a monophyletic group on a phylogenetic tree. These clades may be referred to by
51 a variety of terms, such as ‘subtypes’, ‘genotypes’, ‘groups’, depending on the taxonomic level
52 under investigation or the established scientific literature for the virus in question. The clades
53 usually reflect an attempt to divide pathogen phylogeny and genetic diversity into a set of
54 groupings that are approximately equally divergent, mutually exclusive and statistically well
55 supported. All genome sequences are therefore allocated to one clade or provisionally labelled as

56 'unclassified'. Often multiple hierarchical levels of classification exist for the same pathogens,
57 such as the terms 'type', 'group' and 'subtype' that are used in the field of HIV research.
58

59 Such classification systems are useful for discussing epidemiology and transmission when the
60 number of taxonomic labels remains roughly constant through time; this is the case for slowly-
61 evolving pathogens (for example, many bacteria) and for rapidly-evolving viruses with low rates
62 of lineage turnover (for example, HIV³ and HCV⁴). In contrast, some rapidly-evolving viruses
63 such as influenza A are characterised by high rates of lineage turnover, so that the genetic
64 diversity circulating in any particular year largely emerges out of and replaces the diversity
65 present in the preceding few years. For human seasonal influenza, this behaviour is the result of
66 strong natural selection among competing lineages. In such circumstances a more explicitly
67 phylogenetic classification system is used; for example, avian influenza viruses are classified
68 into 'subtypes', 'clades' and 'higher order clades' according to several quantitative criteria⁵.
69 Such a system can provide a convenient way to refer to the emergence of new (and potentially
70 antigenically-distinct) variants and is suitable for the process of selecting the component viruses
71 for the regularly-updated influenza vaccine. A similar approach to tracking antigenic diversity
72 may be needed to inform SARS-CoV-2 vaccine design efforts. While useful, we recognise that
73 dynamic nomenclature systems based on genetic distance thresholds have the potential to over-
74 accumulate cumbersome lineage names.

75

76 In an ongoing and rapidly changing epidemic, such as SARS-CoV-2, a nomenclature system can
77 facilitate real-time epidemiology by providing commonly-agreed labels to refer to viruses
78 circulating in different parts of the world, thereby revealing the links between outbreaks that

79 share similar virus genomes. Further, a nomenclature system is needed to describe virus lineages
80 that vary in phenotypic or antigenic properties (although it must be stressed that at present there
81 is no conclusive evidence of such variation among currently available SARS-CoV-2 strains).

82

83 **Principles of a dynamic nomenclature system**

84 There are a number of key challenges in the development of a dynamic and utilitarian
85 nomenclature system for SARS-CoV-2. To be valid and broadly accepted a nomenclature needs
86 to: (i) capture local and global patterns of virus genetic diversity in a timely and coherent
87 manner, (ii) track emerging lineages as they move among countries and between populations
88 within each country, (iii) be sufficiently robust and flexible to accommodate new virus diversity
89 as it is generated, and (iv) be dynamic, such that it is able to incorporate both the birth and death
90 of viral lineages through time.

91

92 A special challenge in the case of COVID-19 is that genome sequence data is being generated
93 rapidly and at high volumes, such that by the end of the pandemic we can expect hundreds of
94 thousands of SARS-CoV-2 genomes to have been sequenced. Any lineage naming system must
95 therefore be capable of handling tens to hundreds of thousands of virus genomes sampled
96 longitudinally and densely through time. Further, to be practical, any lineage naming system
97 should have no more than one or two hundred active lineage labels, as any more would obfuscate
98 rather than clarify discussion and will be difficult to conceptualise.

99

100 To fulfil these requirements we propose a workable and practical lineage nomenclature for
101 SARS-CoV-2 that arises from a set of fundamental evolutionary and phylogenetic principles.

102 Some of these principles are, necessarily, specific to the COVID-19 pandemic, reflecting the new
103 reality of large-scale real-time generation of virus genome sequences. The nomenclature system
104 is not intended to represent every evolutionary change in SARS-CoV-2, as these will number
105 many thousand by the end of the pandemic. Instead, the focus is on genetic changes associated
106 with important epidemiological and biological events. Fortunately, because of the early sampling
107 and genome sequencing of COVID-19 cases in China, especially in Hubei province, it appears
108 that the ‘root sequence’ of SARS-CoV-2 is known. Many of the genomes from the earliest
109 sampled cases are genetically identical and hence also likely identical to the most recent common
110 ancestor of all sampled viruses. This occurrence is different to previous viruses and epidemics
111 and provides some advantages for the development of a rational and scalable classification
112 scheme. Specifically, setting the ‘reference sequence’ to be the ‘root sequence’ forms a natural
113 starting point, as direct comparisons in the number and position of mutations can be made with
114 respect to the root sequence.

115

116 During the early phase of the pandemic, it will be possible to unambiguously assign a genome to
117 a lineage through the presence/absence of particular sets of mutations. However, a central
118 component of a useful nomenclature system is that it focuses on those virus lineages that
119 contribute most to global transmission and genetic diversity. Hence, rather than naming every
120 new possible lineage, classification should focus on those that have exhibited onward spread in
121 the population, particularly those that have seeded an epidemic in a new location. For example,
122 the large epidemic in Lombardy, northern Italy, thought to have begun in early February⁶, has
123 since been disseminated to other locations in northern Europe and elsewhere.

124

125 Further, because SARS-CoV-2 genomes are being generated continuously and at a similar pace
126 to changes in virus transmission and epidemic control efforts, we expect to see a continual
127 process of lineage generation and extinction through time. Rather than maintaining a cumulative
128 list of all lineages that have existed since the start of the pandemic, it is more prudent to mark
129 lineages as ‘active’, ‘unobserved’, or ‘inactive’, a designation reflecting our current
130 understanding of whether they are actively transmitting in the population or not. Accordingly,
131 lineages of SARS-CoV-2 documented within the last month are defined here as ‘active’, those
132 last seen >1 month but <3 months ago are classified as ‘unobserved’, and those that have not
133 been seen for >3 months are termed ‘inactive’.

134

135 Although this strategy will allow us to track those lineages that are contributing most to the
136 epidemic, and so reduce the number of names in use, it is important to keep open the possibility
137 that new lineages will appear through the generation of virus genomes from unrepresented
138 locations or from cases with travel history from such locations. For example, the epidemic in
139 Iran, designated B.4 in our system, was identified via returning travellers to other countries⁷.
140 Further, lineages that have not been seen for some time may re-emerge after a period of cryptic
141 transmission in a region. Hence, it is possible for lineages that were previously classified as
142 inactive or unobserved to be later re-labelled as active. We choose the term lineages (rather than
143 ‘clades’, ‘genotypes’ or other designations) for SARS-CoV-2 as it captures the fact that they are
144 dynamic, rather than relying on a static and exclusive hierarchical structure.

145

146 **Lineage naming rules**

147 We propose that major lineage labels begin with a letter. At the root of the phylogeny of SARS-
148 CoV-2 are two lineages that we simply denote as lineages A and B. The earliest lineage A
149 viruses, such as Wuhan/WH04/2020 (EPI_ISL_406801), sampled on 2020-01-05, share two
150 nucleotides (positions 8782 in ORF1ab and 28144 in ORF8) with the closest known bat viruses
151 (RaTG13 and RmYN02). Different nucleotides are present at those sites in viruses assigned to
152 lineage B, of which Wuhan-Hu-1 (GenBank accession MN908947) sampled on 2019-12-26 is an
153 early representative. Hence, although viruses from lineage B happen to have been sequenced and
154 published first⁸⁻¹⁰, it is likely (based on current data) that the most recent common ancestor
155 (MRCA) of the SARS-CoV-2 phylogeny shares the same genome sequence as the early lineage
156 A sequences (e.g. Wuhan/WH04/2020). Importantly, this does *not* imply that the MRCA itself
157 has been sampled and sequenced, but rather that no mutations had accrued between the MRCA
158 and the early lineage A genome sequences. At the time of writing, viruses from both lineages A
159 and B are still circulating in many countries around the world, reflecting the exportation of
160 viruses from Hubei to other regions of China and elsewhere before the strict travel restrictions
161 and quarantine measures were imposed there.

162

163 To add further lineage designations we downloaded 27,767 complete SARS-CoV-2 genomes
164 from the GISAID database¹¹ on 18th May, 2020 and estimated a maximum likelihood tree for
165 these data (see Methods) (Fig. 1). We defined further SARS-CoV-2 lineages, each of which
166 descends from either lineage A or B and is assigned a numerical value (e.g. lineage A.1, or
167 lineage B.2). Lineage designations were made using the following set of conditions:

168

- 169 I. Each descendent lineage should show *phylogenetic evidence* of emergence from an
170 ancestral lineage into another geographically distinct population, implying substantial
171 onward transmission in that population. In the case of a rapidly expanding global lineage
172 the recipient ‘population’ may comprise multiple countries. In the case of large and
173 populous countries it may represent a new region or province. To show *phylogenetic*
174 *evidence* a new lineage must meet *all* of the following criteria: (a) it exhibits one or more
175 shared nucleotide differences from the ancestral lineage, (b) it comprises at least five
176 genomes with >95% of the genome sequenced, (c) genomes within the lineage exhibit at
177 least one shared nucleotide change among them, and (d) a bootstrap value >70% for the
178 lineage defining node. Importantly, criterion (c) helps to focus attention only on lineages
179 with evidence of on-going transmission.
- 180 II. The lineages identified in step (I) can themselves act as ancestors for virus lineages that
181 then emerge in other geographic areas or at later times, provided they satisfy criteria a-d
182 above. This results in a new lineage designation (e.g. A.1.1).
- 183 III. The iterative procedure in step II can proceed for a maximum of three sublevels (e.g.
184 A.1.1.1) after which new descendent lineages are given a letter (in English alphabetical
185 sequence from C, so A.1.1.1.1 would become C.1 and A.1.1.1.2 would become C.2. The
186 rationale for this is that the system is intended only for tracking currently circulating
187 lineages, such that we do not try to capture the entire history of a lineage in its label (that
188 complete history can be obtained by reference to a phylogeny). At the time of writing no
189 C level lineages have been assigned.
- 190 IV. All sequences are assigned to one lineage. For example, if a genome does not meet the
191 criteria for inclusion in a ‘higher level’ lineage (e.g. A.1.2, B.1.3.5) then it is

192 automatically classified into the lowest level for which it does meet the inclusion criteria,
193 which ultimately is ‘A’ or ‘B’.

194

195 Using this scheme we identified 81 viral lineages. These lineages mostly belong to A, B and B.1.

196 We identified six lineages derived from lineage A (denoted A.1-A.6) and two descendant sub-
197 lineages of A.1 (A.1.1 and A.3). We also describe 16 lineages directly derived from lineage B.

198 To date, lineage B.1 is the predominant known global lineage and has been subdivided into > 70
199 sub-lineages. Lineage B.2 currently has six descendant sub-lineages. We are not yet able to
200 further subdivide the other lineages even though some contain very large numbers of genomes.

201 This is because many parts of the world experienced numerous imported cases followed by
202 exponential growth in local transmission. We provide descriptions of these initial lineages,
203 including their geographical locations and time span of sampling, in Table 1. We have also tried
204 to be flexible with the criteria where, for example, the bootstrap value is below 70% but there is
205 strong prior evidence that the lineage exists and is epidemiologically important. In particular, the
206 Italian epidemic comprises two large lineages in our scheme – B.1 and B.2 – reflecting genomes
207 from Italy as well as from large numbers of travellers from these regions and that fall into both
208 lineages.

209

210 A unique and important aspect of our proposed nomenclature is that the status of the currently
211 circulating lineages be assessed at regular intervals, with decisions made about identifying new
212 lineages and flagging those we believe are likely be ‘unobserved’ or ‘inactive’ because none of
213 their members have been sequenced for a considerable time. The names of unobserved or
214 inactive lineages will not be reassigned. These are provisional timescales and the category

215 thresholds may be altered in the future once the dynamics of lineage generation and extinction
216 are better understood. When visualising the epidemic we suggest that these lineages should be no
217 longer labelled to reduce both the number of names in circulation and visual noise, and to focus
218 on the current epidemiological situation.

219

220 **Discussion**

221 While we regard this proposed nomenclature as practical and robust, it is important to recognise
222 that phylogenetic inference carries statistical uncertainty and much of the available genome data
223 is noisy, with incomplete genome coverage and errors arising from the amplification and
224 sequencing processes. We have proposed a genome coverage threshold for proposing new
225 lineages (see above), and we further suggest that sequences are not ascribed a lineage
226 designation unless the genome coverage of that sequence exceeds 70% of the coding region. As
227 noted above, when SARS-CoV-2 genetic diversity is low during the early pandemic period, there
228 will be a direct association between lineage assignment and the presence of particular sets of
229 mutations (with respect to the root sequence). This should help with the development of rapid,
230 algorithmic genome labelling tools. This task will become more complex, but still tractable, as
231 SARS-CoV-2 genetic diversity accumulates, increasing the chance of both homoplasies and
232 reverse mutations. Classification algorithms based on lists of ‘lineage-defining’ mutations may
233 be practical if they are frequently cross-checked and validated against phylogenetic estimations,
234 but will not be as powerful as phylogenetic classification methods that make use of complete
235 genome sequence data to identify relationships. We encourage the research community to
236 develop software and online tools that will enable the automated classification of newly-

237 generated genomes (one such implementation is pangolin, [https://github.com/hCoV-](https://github.com/hCoV-2019/pangolin)
238 [2019/pangolin](https://github.com/hCoV-2019/pangolin)).
239
240 Coronaviruses also frequently recombine, meaning that a single phylogenetic tree may not
241 always adequately capture the evolutionary history of SARS-CoV-2. Although this can make
242 phylogenetic analysis challenging, recombination is readily accommodated within this system of
243 lineage naming and assignment. A distinct recombination event, if it establishes onward
244 transmission, will create a new viral lineage with a distinct common ancestor. Because this new
245 lineage doesn't have a single ancestral lineage they will be assigned the next available
246 alphabetical prefix.

247
248 While we believe that our proposed lineage nomenclature will greatly assist those working with
249 COVID-19, we do not see it as exclusive to other naming systems, particularly those that are
250 specifically intended to track lineages circulating within individual countries for which a finer
251 scale will be helpful. Indeed, there are likely to be strong sampling biases toward particular
252 countries. Further, we note that future genome sequence generation may require adjustments to
253 the current proposal, and any such changes will be detailed at <http://cov-lineages.org/>. We
254 envisage, however, that the general approach described here may be readily adopted for these
255 purposes, and also for other viral epidemics where real-time genomic epidemiology is being
256 undertaken. We expect that this dynamic nomenclature will be most useful for the duration of the
257 global pandemic, which may last a few years. After that time, SARS-CoV-2 will be either
258 globally eliminated or, more likely, become an endemic or seasonal infection. The remaining

259 endemic/seasonal lineages, which will by then be genetically distinct, can simply retain in the
260 post-pandemic period their names from the dynamic nomenclature system.

261

262 **Methods**

263 We downloaded all SARS-CoV-2 genomes (at least 29,000bp in length) from GISAID on May
264 18th 2020. We trimmed the 5' and 3' untranslated regions and retained those genomes with at
265 least 95% coverage of the reference genome (Wuhan-Hu-2019, GenBank accession MN908947).
266 We aligned these sequences using MAFFT's FFT-NS-2 algorithm and default parameter
267 settings¹². We then estimated a maximum likelihood tree using IQ-TREE 2¹³ using the GTR+ Γ
268 model of nucleotide substitution^{14,15}, default heuristic search options, and ultrafast bootstrapping
269 with 1000 replicates¹⁶.

270

271 The maximum likelihood tree and associated sequence metadata were manually curated and the
272 phylogeny was annotated with the lineage designations. This annotated tree, along with a table
273 providing the lineage designation for each genome in the data set, is available for download at
274 <http://cov-lineages.org/>. We also provide a high-resolution PDF figure of the entire tree labelled
275 with lineages. These will be updated on a regular basis. Representative sequences from each
276 lineage were selected to maximise within-lineage diversity and to minimise N-content and used
277 to construct the maximum likelihood tree shown in Figure 1.

278

279 **Correspondence and requests for materials** should be addressed to A.R.

280 (a.rambaut@ed.ac.uk), E.C.H (edward.holmes@sydney.edu.au) and O.G.P.

281 (oliver.pybus@zoo.ox.ac.uk)

282
283 **Acknowledgements.** This work was funded by the Wellcome Trust (Collaboration Award
284 206298/Z/17/Z, ARTIC-Network), the European Research Council under the European Union's
285 Horizon 2020 research and innovation programme (grant agreement no. 725422-
286 ReservoirDOCS) and the European Commission Seventh Framework Programme (FP7/2007-
287 2013)/ European Research Council (grant agreement no. 614725-PATHPHYLODYN), the UK
288 COVID-19 Genomics Consortium (CoG-UK), the Oxford Martin School, and the Australian
289 Research Council (FL170100022).

290
291 **Author Contributions.** A.R., E.C.H. and O.G.P. conceived, designed and supervised the study.
292 A.O'T., J.T.M. and A.R. developed phylogenetic methods. A.R., A.O'T., V.H. J.T.M., C.R., L.
293 du P. and O.G.P. analysed and interpreted the virus genomes. A.R., E.C.H., A.O'T., V.H. and
294 O.G.P. wrote the paper, with input from other authors.

295
296 **Code availability statement.** Details of software and source code that implement the
297 nomenclature system reported here are available at <http://cov-lineages.org>.

298
299 **Data availability statement.** No new data are reported. Virus genome sequences used here are
300 publicly available from <http://gisaid.org>. A table of acknowledgements for the GISAID genome
301 sequences used to develop this work is available at [http://cov-](http://cov-lineages.org/gisaid_acknowledgements.html)
302 [lineages.org/gisaid_acknowledgements.html](http://cov-lineages.org/gisaid_acknowledgements.html)

303
304 **Competing interests.** The authors declare no competing interests.

305

306 **Additional information**

307 **No supplementary Information** is available for this paper.

308

309 **References**

- 310 1. ICTV. Changes to virus taxonomy and the International Code of Virus Classification and
311 Nomenclature ratified by the International Committee on Taxonomy of Viruses (2019).
312 *Arch. Virol.* **2164**, 2417-2429 (2019).
- 313 2. ICTV. ICTV Code: The International Code of Virus Classification and Nomenclature.
314 <https://talk.ictvonline.org/information/w/ictv-information/383/ictv-code/> (2018).
- 315 3. Robertson, D.L. et al. HIV-1 nomenclature proposal. *Science* **288**, 55-56 (2000).
- 316 4. Smith, D.B. et al. Expanded classification of hepatitis C virus into 7 genotypes and 67
317 subtypes: updated criteria and genotype assignment web resource. *Hepatology* **59**, 318–327
318 (2014).
- 319 5. WHO/OIE/FAO H5N1 Evolution Working Group. 2012. Continued evolution of highly
320 pathogenic avian influenza A (H5N1): Updated nomenclature. *Influenza Other Respir.*
321 *Viruses* **6**, 1-5 (2012).
- 322 6. Zehender, G. et al. 2020. Genomic characterisation and phylogenetic analysis of SARS-
323 CoV-2 in Italy. *J. Med. Virol.* doi: 10.1002/jmv.25794. (2020).
- 324 7. Eden, J.-S. et al. An emergent clade of SARS-CoV-2 linked to returned travellers from Iran.
325 *Virus Evol.* **6**, veaa027 (2020).
- 326 8. Lu, R., et al. Genomic characterisation and epidemiology of 2019 novel coronavirus:
327 implications for virus origins and receptor binding. *The Lancet* **395**, 565–574 (2020).
- 328 9. Wu, F. et al. A new coronavirus associated with human respiratory disease in China. *Nature*
329 **579**, 265–269 (2020).
- 330 10. Zhu, N. et al. A novel coronavirus from patients with pneumonia in China, 2019. *New Eng.*
331 *J. Med.* **382**, 727-733 (2020).

- 332 11. Shu, Y. & McCauley, J. GISAID: Global Initiative on Sharing All Influenza Data - from
333 vision to reality. *Eurosurveillance* **22**, 30494 (2017).
- 334 12. Katoh, K., Misawa, K., Kuma, K.-I. & Miyata, T. MAFFT: A novel method for rapid
335 multiple sequence alignment based on fast Fourier transform. *Nuc. Acids Res.* **30**, 3059–66
336 (2002).
- 337 13. Minh, B. Q. et al. IQ-TREE 2: new models and efficient methods for phylogenetic inference
338 in the genomic era. *Mol. Biol. Evol.* **37**, 1530-1534 (2020).
- 339 14. Tavaré, S. Some mathematical questions in biology: DNA sequence analysis. Lectures on
340 mathematics in the life sciences. Vol. 17, 57–86. American Mathematical Society (1986).
- 341 15. Yang, Z. Maximum likelihood phylogenetic estimation from DNA sequences with variable
342 rates over sites: approximate methods. *J. Mol. Evol.* **39**, 306–314. (1994).
- 343 16. Minh, B. Q., Nguyen M. A. T. & von Haeseler, A. Ultrafast approximation for phylogenetic
344 bootstrap. *Mol. Biol. Evol.* **30**, 1188–1195 (2013).
- 345
- 346

347 **Figure Legends**

348

349 **Fig. 1.** Maximum likelihood phylogeny of globally sampled sequences of SARS-CoV-2
350 downloaded from the GISAID database (<http://gisaid.org>) on May 18th 2020. Five representative
351 genomes are included from each of the defined lineages. The largest lineages that are defined by
352 our proposed nomenclature system are highlighted with coloured areas and labelled on the right.
353 The remaining lineages defined by the nomenclature system are denoted by triangles. The scale
354 bar represents the number of nucleotide changes within the coding region of the genome.

355

356 **Table 1.** Proposed nomenclature of early major lineages of SARS-CoV-2. See <https://cov->
 357 [lineages.org/](https://cov-lineages.org/) for full details of each lineage.

358

Lineage	Genomes	Date range	Comments
A	223	Jan-05, Apr-27	Root of the pandemic lies in this lineage, many Chinese sequences with global exports
A.1	1116	Feb-20, Mar-25	Primary outbreak in Washington State, USA
A.2	295	Feb-26, Apr-27	European lineage
A.3	191	Jan-28, Apr-21	USA lineage
A.5	118	Feb-23, Apr-26	European lineage
B	1713	Dec-24, May-03	Base of this lineage lies in China with a lot of global travel between multiple locations
B.1	7438	Jan-24, May-10	Comprises the large Italian outbreak, now represents many European outbreaks, with travel within Europe and from Europe to the rest of the world
B.1.1	6286	Feb-15, May-09	Major European lineage, exports

			to the rest of the world from Europe
B.2	917	Feb-13, May-04	With B.1, comprises the large Italian outbreak
B.3	752	Feb-23, Apr-23	UK lineage
B.4	258	Jan-18, Apr-14	Likely the primary Iranian outbreak

359

360

