1 # Rapid whole genome sequence typing reveals multiple waves of

2 # SARS-CoV-2 spread

3

4 Ahmed M. Moustafa[1], Paul J. Planet[1,2,3*]

5

6 **1. Division of Pediatric Infectious Diseases, Children's Hospital of Philadelphia,**

7 **Philadelphia, PA 19104, USA.**

8

9 **2. Department of Pediatrics, Perelman College of Medicine, University of**

10 **Pennsylvania, Philadelphia, PA 19104, USA.**

11

12 **3. Sackler Institute for Comparative Genomics, American Museum of Natural**

13 **History, New York, NY 10024, USA.**

14

15 **Emails**

16 **AMM: moustafaam@email.chop.edu**

17 **PJP: planetp@email.chop.edu**

18

19 **\*Corresponding Author**

20

**Abstract**

As the pandemic SARS-CoV-2 virus has spread globally its genome has diversified to

an extent that distinct clones can now be recognized, tracked, and traced. Identifying

clonal groups allows for assessment of geographic spread, transmission events, and

identification of new or emerging strains that may be more virulent or more

transmissible. Here we present a rapid, whole genome, allele-based method (GNUVID)

for assigning sequence types to sequenced isolates of SARS-CoV-2 sequences. This

sequence typing scheme can be updated with new genomic information extremely

rapidly, making our technique continually adaptable as databases grow. We show that

our method is consistent with phylogeny and recovers waves of expansion and

replacement of sequence types/clonal complexes in different geographical locations.

GNUVID is available as a command line application

(https://github.com/ahmedmagds/GNUVID).


**Keywords**

*SARS-CoV-2, COVID-19, nomenclature, lineages, WhatsGNU, cgMLST, wgMLST,*

*clonal complex*


**Introduction**

Rapid sequencing of the SARS-CoV-2 pandemic virus has presented an

unprecedented opportunity to track the evolution of the virus and to understand the

emergence of a new pathogen in near-real time. During its explosive radiation and

global spread, the virus has accumulated enough genomic diversity that we are now

2

44    able to identify distinct lineages and track their spread in distinct geographic locations

45    and over time [1-6]. Phylogenetic analyses in combination with rapidly growing

46    databases [1, 7] have been instrumental in identifying distinct clades and tracing how

47    they have spread across the globe, as well as estimating calendar dates for the

48    emergence of certain clades [1-4]. This information is extremely useful in assessing the

49    impact of early measures to combat spread as well as identifying missed opportunities

50    [3]. Going forward whole genome sequences will be useful for identifying emerging

51    clones or hotspots of reemergence.

52        In all of these efforts, identification of specific clones, clades, or lineages, is a

53    critical first step, and there are few systems available to do this [1]. As of June $1_{st}$ there

54    are already 35,291 and 4,636 complete genomes (>29,000bp) available at GISAID [7]

55    and GenBank [8], respectively. To address the problem of identifying sequence types in

56    SARS-CoV-2 and leverage these huge datasets, we took inspiration from a an

57    approach used widely in bacterial nomenclature, multilocus sequence typing (MLST) [9].

58    Our panallelome approach to developing a whole genome (wgMLST) scheme for

59    SARS-CoV-2 uses a modified version of our recently developed tool, WhatsGNU [10],

60    to rapidly assign an allele number to each gene nucleotide sequence in the virus's

61    genome creating a sequence type (ST). The ST is codified as the sequence of allele

62    numbers for each of the 10 genes in the viral genome.

63        Here we show that this approach allows us to link STs into clearly defined clonal

64    complexes (CC) that are consistent with phylogeny. We show that assessment of STs

65    and CCs agrees with multiple introductions of the virus in certain geographical locations.

66    In addition, we use temporal assessment of STs/CCs to uncover waves of expansion

67    and decline, and the apparent replacement of certain STs with emerging lineages in

68    specific geographical locations.

69

70    **Results and Discussion**

71         We developed the GNU-based Virus IDentification (GNUVID) system as a tool

72    that automatically assigns a number to each unique allele of the 10 open reading

73    frames (ORFs) of SARS-CoV-2 (**Figure 1A**) by modifying our tool WhatsGNU [10]

74    (**Supplementary Methods**). GNUVID compressed the 104,220 ORFs in 10,422 high

75    quality GISAID genomes (**Supplementary Table 1**) to 6244 unique alleles in less than

76    one minute on a standard desktop, achieving 17-fold compression and losing no

77    information. The majority of these alleles (65%) are for ORF1ab which represents 71%

78    of the genome length (**Figure 1A**). Strikingly, the most abundant alleles of each ORF

79    (except ORF1ab) were present in at least 79% of the 10,422 isolates, and for 8 ORFs

80    (ORF3a-10) the allele that was observed in the earliest genomes was also the most

81    prevalent, suggesting strong nucleotide level conservation over time.

82         Some widespread alleles corresponded to mutations that have been

83    hypothesized to be important to the evolution or pathogenesis of the virus. For instance,

84    for the S gene, the gene for the Spike protein, 64% (526/817) of unique alleles have the

85    A23403G (D614G) mutation (**Figure 1B**) that has been associated with the emergence

86    of increased transmission whether through increased transmissibility [11] or lapses in

87    control around this variant [3]. The first allele isolated and sequenced (allele 17) that

88    carries this mutation was first recorded on January 24th in China. The most common S

89    gene allele that carries the A23403G mutation (allele 26) was present in 55% of the

4

90    isolates. For ORF3a, which was shown to activate the NLRP3 inflammasome [12], 35%

91    (126/357) of alleles have the G25563T (Q57H) mutation representing 33% of the

92    isolates. The earliest sequenced, and most common, ORF3a allele that carries this

93    mutation (allele 25) was isolated in France on February 21st in a virus that also carries

94    also the A23403G mutation in the spike gene.

95       To create an ST for each isolate GNUVID automatically assigned 5510 unique

96    ST numbers based on their allelic profile (**Supplementary Table 2**). We then used a

97    minimum spanning tree (MST) to group STs into larger taxonomic units, clonal

98    complexes (CCs), which we define here as clusters of >20 STs that are single or double

99    allele variants away from a "founder". Using the goeBURST algorithm [13, 14] to build

100   the MST and identify founders, we found 24 CCs representing 79% (4352/5510) of all

101   unique STs (**Figure 1B**).

102      When the global region of origin for each genome sequence was mapped to

103   each CC there was a strong association of some CCs with certain geographical

104   locations. For instance, genomes from CCs 255, 300, 301, 317, 348, 355, 369, 399,

105   454, 498, 985, 1063, 1148 are predominantly from Europe while genomes from CCs 26,

106   800 and 927 are mainly from Asia (**Figure 1B**). Interestingly, genomes originating from

107   the US appear to be associated with 2 very divergent CCs, potentially reflecting two

108   major introductions. The first, CC256, is associated with locations on the West Coast,

109   specifically Washington state. The first two isolates belonging to CC256 are from China

110   followed by the first isolate from Washington (01/19/2020). The second predominant US

111   CC, CC258, is closely related to other CCs found predominantly in Europe (**Figure 1B**

112   **and 1C**). Isolates of CC258 were initially found and sequenced in Europe, followed by

5

113    the US East Coast, and later in other US locations (**Figure 1B**). Interestingly, almost all

114    isolates (99%) from CC258 and its descendants CCs 768, 800, 844 and 1063 (**Figure**

115    **1B**) carry the G25563T mutation in ORF3a, representing 88% of all isolates that carry

116    this mutation; the other 12% are from STs that were not assigned CCs. CC800 is

117    interesting for its geographic predominance in the Middle East (75% from Saudi Arabia

118    and Turkey) and its close relationship to ST338 and ST258, which are mostly found in

119    the US. This may signal a transmission event from the US to the Middle East.

120          To show that CCs are mostly consistent with whole genome phylogenetic trees

121    we produced a maximum likelihood tree and mapped the CC designations onto the tree.

122    Figure 1C shows that members of the same CC usually grouped together in clades

123    (**Supplementary File 4**). One limitation of any ST/CC classification strategy is that

124    paraphyletic groups can occur as a new ST arises from an older ST (e.g. CC301

125    emerged from CC255 making CC255 paraphyletic). While this means not all ST/CC

126    groups will be monophyletic, this property of the nomenclature may be helpful in

127    gauging emergence and replacement of an ancestral form.

128          To further validate our wgMLST classification system we compared it to the

129    recently proposed "dynamic lineages nomenclature" for SARS-CoV-2 using the pangolin

130    application[1]. A high percentage of viruses (90.5%;40-100%) with the same CC were

131    assigned to the same lineage. When sublineages of the dominant lineage designation

132    were included, this average rose to 99% (89-100%), showing strong agreement

133    between these classification schemes (**Supplementary Table 2**).

134          Because we included collection dates for each genomic sequence, we can use

135    STs and CCs to better understand the emergence and replacement of certain lineages

136    in certain geographical regions over time. **Figure 2A** shows temporal plots of the most

137    common 12 CCs around the world. This makes clear the emergence of new CCs over

138    time such as CC255, CC300 and CC258. CC4, the earliest CC, started by representing

139    60% of sequenced genomes in mid-January, but had dropped to only 5% by mid-March.

140         Of course, relative proportions of STs or CCs isolated and sequenced may be a

141    highly biased statistic that is contingent upon where the isolate comes from, the

142    decision to sequence its genome, and the local capacity to sequence a whole genome.

143    Certain regions (US and Europe) clearly sequenced more genomes later in the

144    pandemic compared to other countries.

145         Focusing on specific geographic regions may help to partially ameliorate this

146    bias, and we chose to focus on three different regions (China, Europe and the US). The

147    temporal plot of China shows expansion of local clones (CC4 and CC256) that likely

148    spread to other countries early in the pandemic and then decreased in China over time.

149    In contrast, two new CCs 927 and 454 appear to have emerged more recently with

150    earliest isolation dates of March 18th and April 16th, respectively, though this should be

151    interpreted with caution because few sequences (n=7 and 6) were available/included.

152    Interestingly, CC258 was first isolated in China in mid-March while it already

153    represented 14% of the genomes in Europe by the end of February (**Figure 2B and D**),

154    potentially reflecting transmission of new lineages back to China later in the pandemic.

155    By the end of January, although CC4 represented 39% of the sequenced genomes in

156    China, only one isolate (1/6) of CC4 was isolated in Wuhan, showing different patterns

157    of circulating clones at the same timepoint in different parts of the same country (**Figure**

158    **2B and C**).

159      Interestingly, Europe showed a general CC diversity over time resembling that of

160    the worldwide temporal plot, and then showed expansion of the local CC300 and

161    CC255 after mid-February (**Figure 2D**).

162      The US plot (**Figure 2E**) reflects the two possible introductions on the west and

163    east coasts from Asia and Europe, respectively, with the current dominance (more than

164    45%) of CC258. Focusing on Washington, it is interesting to note the possible

165    replacement of CC256 by CC258 perhaps by introduction from the East Coast or

166    Europe (**Figure 2F**) [2, 4]. In New York, a different pattern is seen with CC258 being

167    persistently dominant (**Figure 2G**). However, a more granular view of STs in New York,

168    not CCs, shows a shifting epidemiology with ST258 declining and the rise of closely

169    related SLVs and DLVs of ST258 (**Figure 2H**).

170      While our wgMLST approach is rapid and robust it has several limitations.

171    Because a change in any allele creates a new ST our method may accumulate and

172    count "unnecessary" STs that have been seen only once or may be due to a

173    sequencing error. This is partially ameliorated by the use of the CC definition that allows

174    some variability amongst the members of a group. A large number of STs also may

175    allow more granular approaches to tracking new lineages. Our method is also limited by

176    the quality and extent of the database. For this implementation we limited the database

177    to genomes that do not have any ambiguity or degenerate bases. However, these

178    genomes could be queried through our tool to be assigned to the closest ST/CC.

179    Another limitation is the stability of the classification system, some virus genomes may

180    be reassigned to new CCs as clones expand epidemiologically, but this may also reflect

181    a dynamic strength as circulating viruses emerge and replace older lineages.

182

**Conclusion**

183

184     The genomic epidemiology of the 10,422 SARS-CoV-2 isolates studied here

185     show six predominant CCs circulated/circulating globally. Our tool (GNUVID) allows for

186     fast sequence typing and clustering of whole genome sequences in a rapidly changing

187     pandemic. As illustrated above, this can be used to temporally track emerging clones or

188     identify the likely origin of viruses. With stored metadata for each sequence on date of

189     isolation, geography, and clinical presentation, new genomes could be matched almost

190     instantaneously to their likely origins and potentially related clinical outcomes.

191

192     **Methods**

193     All SARS-CoV-2 genomes (n=17,504) that are complete and high coverage were

194     downloaded from GISAID [7] on May 17th 2020. We kept 16,866 that were at least

195     29,000 bp in length and had less than 1% "N"s. Our wgMLST scheme was composed of

196     all 10 ORFs in the SARS-CoV-2 genome [15]. The 10 ORFs were identified in the

197     remaining 16,866 genomes using blastn [16] and any genome that had any ambiguity or

198     degenerate bases (any base other than A,T,G and C) in the 10 open reading frames

199     (ORF) was excluded. The remaining 10,422 genomes were fed to the GNUVID tool in a

200     time order queue (first-collected to last-collected), which assigned a ST profile to each

201     genome. The identified STs by GNUVID were fed into the PHYLOViZ tool [17] to identify

202     CCs at the double locus variant (DLV) level using the goeBURST MST [13, 14]. CCs

203     were mapped back to the STs using a custom script. Pie charts were plotted using a

204    custom script. Temporal plots were extracted using a custom script and plotted in

205    GraphPad Prism v7.0a.

206

207    To show the relationship between our typing scheme and phylogeny, we constructed a

208    maximum likelihood tree. Briefly, we masked the 5' and 3' untranslated regions in the

209    10,422 genomes. We aligned these sequences using MAFFT's FFT-NS-2 algorithm

210    (options: --add --keeplength) [18] to the reference MN908947.3 [15]. A maximum

211    likelihood tree using IQ-TREE 2 [19] was estimated using the HKY model of nucleotide

212    substitution [20], default heuristic search options, and ultrafast bootstrapping with 1000

213    replicates [21]. The tree was rooted to MN908947.3. The tree and ST/CC data were

214    visualized in iTOL [22]. We assigned a lineage [1] to each member of the 24 CCs using

215    pangolin (https://github.com/hCoV-2019/pangolin) (options: -t 8). The GNUVID

216    database will be updated weekly with new added high-quality genomes from GISAID

217    [7]. Detailed methods are in **Additional file 1**.

218

219    **List of abbreviations**

220    WhatsGNU    What is Gene Novelty Unit

221    GNUVID        Gene Novelty Unit-based Virus Identification

222    ST      Sequence Type

223    CC      Clonal Complex

224    SARS-CoV-2 Severe Acute Respiratory Syndrome Corona Virus 2

225    COVID-19    Corona Virus Disease 2019

226    MLST Multilocus Sequence Typing

227    cgMLST core genome MLST

228    wgMLST whole genome MLST

229

230    **Figure Legends**

231    **Figure 1.** Sequence Typing Scheme for SARS-CoV-2. **A.** Map of SARS-CoV-2 virus

232    genome showing the length in base pairs (bp) of the 10 ORFs, numbers of alleles in the

233    current database, and the prevalence of the top two alleles of each ORF in the 10,422

234    database isolates. **B.** Minimum spanning tree from goeBURST of the 5510 Sequence

235    Types (STs) showing the 24 Clonal Complexes (CCs) identified in the dataset. The

236    largest six CCs are red and the other 18 CCs are in black. The pie charts show the

237    percentage distribution of genomes from the different geographic regions in each CC.

238    The letter A and G next to the pie charts represent the Spike ORF nucleotide at position

239    23403 in MN908947.3. The ancestral nucleotide is A and the mutation is G resulting in

240    D614G amino acid change. At least 98% of the genomes of each CC had the reported

241    nucleotide (except for CC26 where it was 93%). **C.** Maximum likelihood phylogeny of

242    the 10,422 global high-quality SARS-CoV-2 sequences downloaded from the GISAID

243    database (http://gisaid.org) on May 17th 2020 (**Supplementary Table 1**). The tree is

244    rooted on the reference sequence MN908947.3. The tree was visualized in iTOL. Only

245    the most common seven CCs were shown for easier visualization. Nodes with 200-500

246    leaves were collapsed for better visualization. The raw tree is available as

247    **Supplementary File 4**.

248

11

249    **Figure 2.** Temporal Plots of circulating STs/CCs at different geographical locations

250    (Global, China, Wuhan, Europe, USA, Washington, NY (CC) and NY (ST)). The

251    visualizations were limited to the most common CCs/STs. The raw data can be obtained

252    from the authors upon request.

253

254    **Additional files**

255    **Additional file 1: Supplementary Methods** (txt, 34 Kb)**.**

256    **Additional file 2: Table S1. Acknowledgment Table** (xls, 2.1 Mb)**.**

257    **Additional file 3: Table S2. GNUVID Database Strains Report Table** (xlsx, 778 Kb).

258    **Additional file 4: Maximum Likelihood Tree of the 10422 strains** (nex, 369 Kb).

259

260    **Availability of data and material**

261    The compressed genomes from our quality controlled dataset are available from the

262    corresponding author and available online for download. The compressed database will

263    be updated weekly on https://github.com/ahmedmagds/GNUVID. Source code for

264    GNUVID can be found in its most up-to-date version here,

265    https://github.com/ahmedmagds/GNUVID, under the GNU General Public License.

266

267    **Competing interests**

268    The authors declare that they have no competing interests

269

270    **Funding**

271   PJP and AMM are supported by NIH 1R01AI137526-01, PLANET19G0,

272   1R21AI144561-01A1. PJP is further supported by R01NR015639,

273

274   **Authors' contributions**

275   Conceptualization: AMM, PJP; Coding: AMM; Writing – Reviewing and Editing: AMM,

276   PJP.

277

278   **Acknowledgements**

283

284   **References**

285   1.   Rambaut A, Holmes EC, Hill V, O'Toole Á, McCrone JT, Ruis C, du Plessis L, Pybus OG: **A**
286        **dynamic nomenclature proposal for SARS-CoV-2 to assist genomic epidemiology.**
287        *bioRxiv* 2020:**2020.2004.2017.046086.
288   2.   Deng X, Gu W, Federman S, Du Plessis L, Pybus O, Faria N, Wang C, Yu G, Pan C-Y,
289        Guevara H, et al: **A Genomic Survey of SARS-CoV-2 Reveals Multiple Introductions into**
290        **Northern California without a Predominant Lineage.** *medRxiv*
291        2020:**2020.2003.2027.20044925.
292   3.   Worobey M, Pekar J, Larsen BB, Nelson MI, Hill V, Joy JB, Rambaut A, Suchard MA,
293        Wertheim JO, Lemey P: **The emergence of SARS-CoV-2 in Europe and the US.** *bioRxiv*
294        2020:**2020.2005.2021.109322.
295   4.   Bedford T, Greninger AL, Roychoudhury P, Starita LM, Famulare M, Huang M-L, Nalla A,
296        Pepper G, Reinhardt A, Xie H, et al: **Cryptic transmission of SARS-CoV-2 in Washington**
297        **State.** *medRxiv* 2020:**2020.2004.2002.20051417.
298   5.   Shen L, Dien Bard J, Biegel JA, Judkins AR, Gai X: **Comprehensive genome analysis of**
299        **6,000 USA SARS-CoV-2 isolates reveals haplotype signatures and localized**
300        **transmission patterns by state and by country.** *medRxiv*
301        2020:**2020.2005.2023.20110452.

302    6.    Chen Z-w, Li Z, Li H, Ren H, Hu P: **Global genetic diversity patterns and transmissions of**
303          **SARS-CoV-2.** *medRxiv* 2020:2020.2005.2005.20091413.
304    7.    Shu Y, McCauley J: **GISAID: Global initiative on sharing all influenza data - from vision**
305          **to reality.** *Euro Surveill* 2017, **22**.
306    8.    Sayers EW, Cavanaugh M, Clark K, Ostell J, Pruitt KD, Karsch-Mizrachi I: **GenBank.**
307          *Nucleic Acids Res* 2019, **47**:D94-D99.
308    9.    Maiden MC, Bygraves JA, Feil E, Morelli G, Russell JE, Urwin R, Zhang Q, Zhou J, Zurth K,
309          Caugant DA, et al: **Multilocus sequence typing: A portable approach to the**
310          **identification of clones within populations of pathogenic microorganisms.** *PNAS* 1998,
311          **95**:3140-3145.
312    10.   Moustafa AM, Planet PJ: **WhatsGNU: a tool for identifying proteomic novelty.** *Genome*
313          *Biology* 2020, **21**:58.
314    11.   Korber B, Fischer WM, Gnanakaran S, Yoon H, Theiler J, Abfalterer W, Foley B, Giorgi EE,
315          Bhattacharya T, Parker MD, et al: **Spike mutation pipeline reveals the emergence of a**
316          **more transmissible form of SARS-CoV-2.** *bioRxiv* 2020:2020.2004.2029.069054.
317    12.   Siu K-L, Yuen K-S, Castaño-Rodriguez C, Ye Z-W, Yeung M-L, Fung S-Y, Yuan S, Chan C-P,
318          Yuen K-Y, Enjuanes L, Jin D-Y: **Severe acute respiratory syndrome coronavirus ORF3a**
319          **protein activates the NLRP3 inflammasome by promoting TRAF3-dependent**
320          **ubiquitination of ASC.** *The FASEB Journal* 2019, **33**:8865-8877.
321    13.   Francisco AP, Bugalho M, Ramirez M, Carriço JA: **Global optimal eBURST analysis of**
322          **multilocus typing data using a graphic matroid approach.** *BMC Bioinformatics* 2009,
323          **10**:152.
324    14.   Feil EJ, Li BC, Aanensen DM, Hanage WP, Spratt BG: **eBURST: Inferring Patterns of**
325          **Evolutionary Descent among Clusters of Related Bacterial Genotypes from Multilocus**
326          **Sequence Typing Data.** *Journal of Bacteriology* 2004, **186**:1518.
327    15.   Wu F, Zhao S, Yu B, Chen YM, Wang W, Song ZG, Hu Y, Tao ZW, Tian JH, Pei YY, et al: **A**
328          **new coronavirus associated with human respiratory disease in China.** *Nature* 2020,
329          **579**:265-269.
330    16.   Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ: **Basic local alignment search tool.**
331          *Journal of Molecular Biology* 1990, **215**:403-410.
332    17.   Nascimento M, Sousa A, Ramirez M, Francisco AP, Carrico JA, Vaz C: **PHYLOViZ 2.0:**
333          **providing scalable data integration and visualization for multiple phylogenetic**
334          **inference methods.** *Bioinformatics* 2017, **33**:128-129.
335    18.   Katoh K, Misawa K, Kuma K, Miyata T: **MAFFT: a novel method for rapid multiple**
336          **sequence alignment based on fast Fourier transform.** *Nucleic Acids Res* 2002, **30**:3059-
337          3066.
338    19.   Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, von Haeseler A,
339          Lanfear R: **IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in**
340          **the Genomic Era.** *Mol Biol Evol* 2020, **37**:1530-1534.
341    20.   Hasegawa M, Kishino H, Yano T: **Dating of the human-ape splitting by a molecular clock**
342          **of mitochondrial DNA.** *J Mol Evol* 1985, **22**:160-174.
343    21.   Hoang DT, Chernomor O, von Haeseler A, Minh BQ, Vinh LS: **UFBoot2: Improving the**
344          **Ultrafast Bootstrap Approximation.** *Mol Biol Evol* 2018, **35**:518-522.

345    22.    Letunic I, Bork P: **Interactive Tree Of Life (iTOL) v4: recent updates and new**
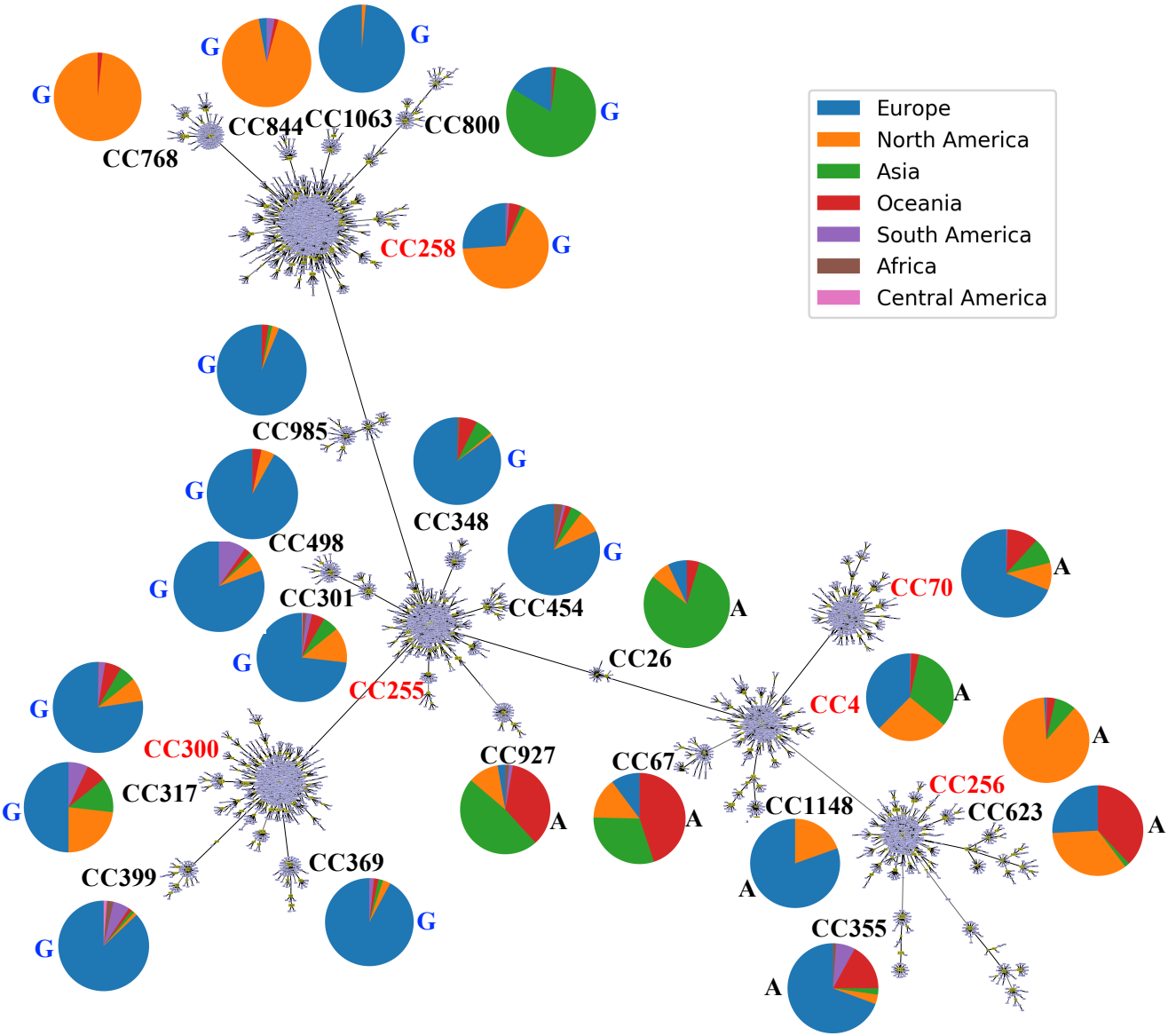346           **developments.** *Nucleic Acids Res* 2019, **47:**W256-W259.
347

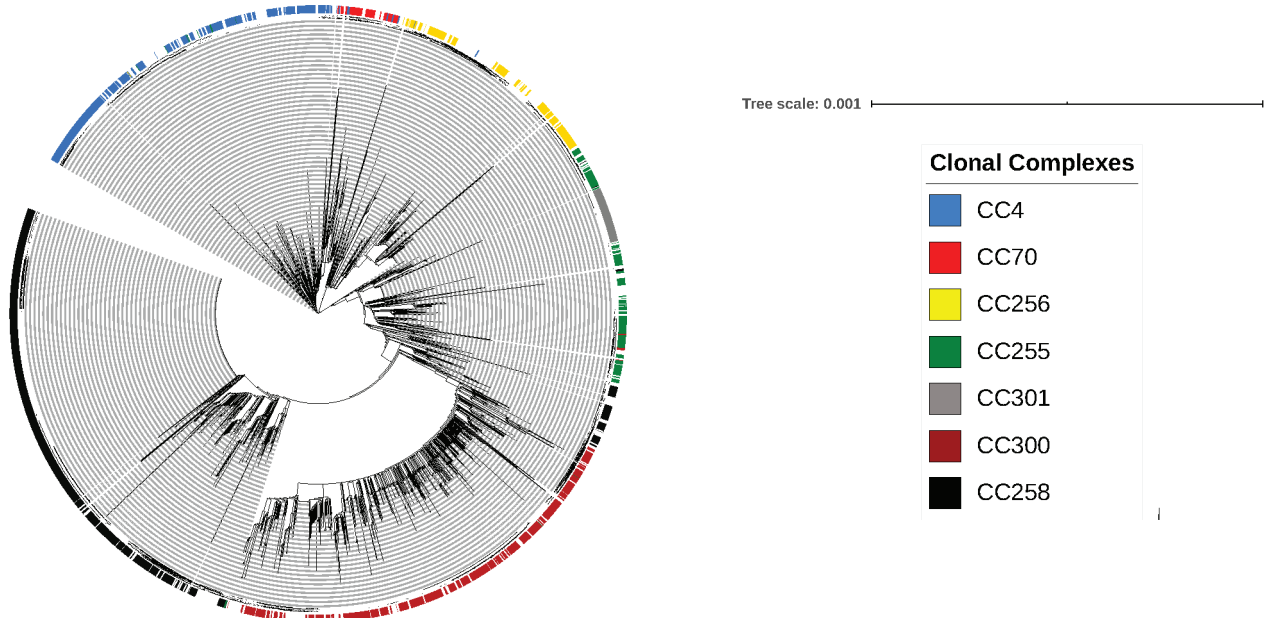# Figure 1 Sequence Typing Scheme for SARS-CoV-2

**A**

|  | 1ab | S | 3a | E | M | 6 | 7a | 8 | N | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| Length (bp) | 21290 | 3822 | 828 | 228 | 669 | 186 | 366 | 366 | 1260 | 117 |
| Number of Alleles | 4050 | 817 | 357 | 50 | 139 | 50 | 116 | 127 | 495 | 43 |
| Top Two Alleles Prevalence (%) | 19.5 | 79 | 83.6 | 98.9 | 92.6 | 98.7 | 97.6 | 93.7 | 82.1 | 98.9 |

**B**



**C**