

1 **Population genomic insights into the evolution of the SARS-CoV-2**

2 **Omicron variant**

3

4 Kritika M. Garg^{1,2}, Vinita Lamba³, Balaji Chattopadhyay^{2,3,*}

5 *¹Centre for Interdisciplinary Archaeological Research, Ashoka University, Sonipat, Haryana,*
6 *131029, India*

7 *²Department of Biology, Ashoka University, Sonipat, Haryana, 131029, India*

8 *³Trivedi School of Biosciences, Ashoka University, Sonipat, Haryana, 131029, India*

9

10

11 *Corresponding author: Balaji Chattopadhyay

12 Email ID: balaji.chattopadhyay@ashoka.edu.in

13 ORCID: 0000-0002-4423-3127

14

15

16

17

18

19

20

21

22

23

24

25

26 **ABSTRACT**

27 A thorough understanding of the patterns of population subdivision of a pathogen can prevent
28 disease spread. For SARS-CoV-2, the availability of millions of genomes makes this task
29 analytically challenging. Our study used population genomic methods and identified subtle
30 subdivisions within the Omicron variant, in addition to that captured by the Pango lineage.
31 Further, some of the identified clusters of the Omicron variant revealed statistically
32 significant signatures of selection or expansion revealing the role of microevolutionary
33 processes in the spread of the virus. These are crucial information for policy makers as
34 preventive measures can be designed to mitigate further spread based on a holistic
35 understanding of the variability of the virus and evolutionary processes aiding its spread.

36

37 **Running title:** Genetic subdivision of Omicron lineage

38

39 **KEY WORDS:** haplotype network, population subdivision, SARS-CoV-2, Omicron

40

41

42

43

44

45

46

47

48

49

50

51 **MAIN TEXT**

52

53 The exponential increase in the number of SARS-CoV-2 genomes has brought with it a
54 unique set of challenges for data analysis. With over ten million genome sequences already
55 available, new algorithms are being designed to tackle the deluge of data. Most available
56 analytical tools are designed to provide the overall evolutionary relationship between various
57 lineages. However, obtaining a finer level understanding of the diversity and subdivision
58 within a lineage can provide important insights into pathogen evolution, particularly during
59 ongoing pandemics [1]. In this study, we utilized population genomic methods to understand
60 the subdivision of the Omicron lineage of SARS-CoV-2 across the globe in an attempt to
61 elucidate the evolutionary history of the variant.

62

63 We downloaded all SARS-CoV-2 genome sequences belonging to the Omicron lineage
64 available up to 31 January 2022 from the GISAID repository (<https://www.gisaid.org/>, see
65 appendix Table S1) and cleaned them using Nextclade CLI [2] (supporting information). We
66 retained 14,002 good quality sequences, most of which were from Denmark (n=11,272). We
67 aligned the filtered SARS-CoV-2 genomes to the Wuhan reference genome (accession ID:
68 MN908947.1) in Nextalign CLI [2] using default parameters. Further, we assigned the
69 lineage for each sequence using the pangolin web server version 3.1.20 and 4.0.6 accessed on
70 3rd March and 6th May 2022 respectively [1].

71

72 We used two different approaches to understand the population subdivision within our
73 SARS-CoV-2 dataset. For the first approach, we constructed a haplotype network using the
74 program VENAS (Viral genome Evolution Network Analysis System) to generate the SARS-
75 CoV-2 evolution and transmission network [3]. VENAS can analyze thousands of genomes

76 collected in a short span of time (in few minutes) and is a useful tool for tracking changes
77 across a transmission chain. It identifies mutations across alignments and constructs a
78 network based on hamming distances. In VENAS we first estimated the effective parsimony-
79 informative sites (ePIS) and minor allele frequencies (MAF) using default settings and
80 retained 5,253 genomes. These were then used to construct an evolutionary network which
81 was viewed in Cytoscape 3.9.1 using perforce-force directed method [4]. Default filtering
82 settings in VENAS removed all sequences belonging to BA.1 lineage and hence, we analyzed
83 sequences from this lineage separately (n = 260 genomes for subdivision analysis).

84

85 For the second approach we used the discriminant analysis of principal components (DAPC)
86 method to understand the fine-scale subdivision patterns observed in each lineage (based on
87 Pango lineage definitions mentioned above). DAPC is a useful method to detect subdivision
88 as it maximizes between-group differences while minimizing within-group variability [5]. It
89 is a relatively fast method to detect complex subdivision patterns from genomic data. We
90 used the filtered genomes obtained from VENAS and performed DAPC analysis for both
91 lineages using R adegenet package [5,6]. We first identified the optimal number of clusters
92 within each dataset using K-means algorithm and then performed DAPC analysis. We further
93 identified the unique mutations for each of the DAPC clusters and only considered mutation
94 which were present in at least 70% of the sequences belonging to the cluster.

95

96 To estimate the level of genomic diversity within our dataset, we characterized all
97 substitutions in reference to the Wuhan genome using VENUS. We considered all 14,003
98 good quality sequences and identified the mutations for the 12 functional open reading
99 frames (ORF). We further estimated Tajima's D for the spike protein sequences for all the
100 clusters identified in DAPC using MEGA 10.2.6 [7]. Tajima's D test is widely used to

101 identify signatures of microevolutionary forces like population fluctuations and selection
102 acting upon populations. We estimated Tajima's D for each genetic cluster identified by
103 DAPC to avoid confounding effects of population subdivision.

104

105 We observed similar pattern of subdivision for each Omicron lineage using both VENUS and
106 DAPC. While VENAS produced numerous nodes/groups (111 for BA.1 and 1,046 BA.2),
107 these were nested within fewer broader subdivisions retrieved by DAPC (Fig. 1). We
108 identified five clusters for the BA.1 lineage and ten clusters for BA.2 lineage. However,
109 when we visualized our results, we observed that two clusters for the BA.1 lineage were
110 clubbed together and four clusters for the BA.2 lineage were clubbed together (Fig. 1B and
111 1D). Further, no sequences were assigned to clusters 1 and 3 for the BA. 2 lineage. Thus,
112 effectively only four clusters for both BA.1 and BA.2 lineages were identified by the DAPC
113 method. Most of these clusters shared mutations and only a few private mutations were
114 identified (Fig. 1B and 1D). Although, both the methods use different approaches, together
115 they provide a robust understanding of the finer subdivision patterns within fast evolving
116 lineages.

117

118 At the start of the study, Pango lineage definition 3.1.20 was available which divided the
119 sequences into BA.1, BA.1.1, and BA.2 lineages and our population genomic based
120 clustering identified fine-scale subdivision within these lineages. With the updated Pango
121 lineages version 4.0.6 which is available now, there is broad agreement with the lineage
122 definitions between our methodology (DAPC+ VENAS) and Pango lineage. Using our
123 population genomics pipeline, we could identify mutations used for Pango lineage
124 definitions. For example, G22599A and G5924A are defining mutations for BA.1.1 (cluster 1
125 and 5) and BA.1.17 (cluster 4) respectively (Fig. 1B). Our method also identified fine-scale

126 subdivision within the Pango definitions, for example, DAPC clusters 5, 7, and 10 all
127 correspond to BA.2.9 lineage from Denmark (Fig. 1D) and were segregated based on three
128 unique mutations (Fig. 1D). One of the mutations unique to cluster 7, C22570T is also used
129 to define the BA.2.9 lineage. However, within our dataset this mutation was only present in
130 cluster 7, highlighting finer subdivision within this lineage (Fig. 1D).

131

132 A detailed inspection of the pattern of substitution among the study genomes revealed that as
133 expected, the spike protein, ORF1a, and ORF1b harbored the maximum number of variations
134 (Fig. S1). The most frequent mutations observed were C to T transition and G to T
135 transversion (Fig. S2). Tajima's D values for the spike protein sequences were negative for
136 all the DAPC clusters, with significant values observed only for four clusters (Table S2).

137 Significant negative Tajima's D suggests that these DAPC clusters (BA.1 lineage: cluster 1;
138 BA.2 lineage: cluster 2, 5, and 6) are undergoing rapid expansions from a small population
139 size and/or have experienced recent selective sweeps, making them potential target for
140 surveillance and monitoring programs. Interestingly, cluster 2 of BA.1 lineage which did not
141 exhibit a signature of expansion and/or selection, also harbors three unique mutations
142 (C25708T, A29301G, T10135C), which have been identified as suppressors of the spread of
143 this variant by Yang et al. [8] (Fig. 1A and 1B; Table S2). In conjunction, these results
144 identified evolutionary processes affecting virus transmission, highlighting the importance of
145 our approach in identifying sub-lineages of concern and hotspots of spread of SARS-CoV-2.
146 Using a combination of population genomic methods, we could recover subtle variations
147 (within established lineage definitions), some of which can spread more than others. This
148 provides an easy analytical framework which can be used by policy makers to identify
149 variants of potential concern thereby facilitating disease mitigation.

150

151 **ACKNOWLEDGEMENTS**

152 B.C. acknowledges the startup funding from Trivedi School of Biosciences (TSB), Ashoka
153 University, India and K.M.G. acknowledges the support from the DBT-Ramalingaswami
154 Fellowship (No. BT/HRD/35/02/2006). V.L. was supported by TSB fellowship. We
155 gratefully acknowledge the following Authors from the Originating laboratories responsible
156 for obtaining the specimens and the Submitting laboratories where genetic sequence data
157 were generated and shared via GISAID Initiative, on which this research is based. A full
158 acknowledgement table can be found in Appendix Table S1.

159

160 **DECLARATION OF INTEREST STATEMENT**

161

162 The authors declare no conflict of interest

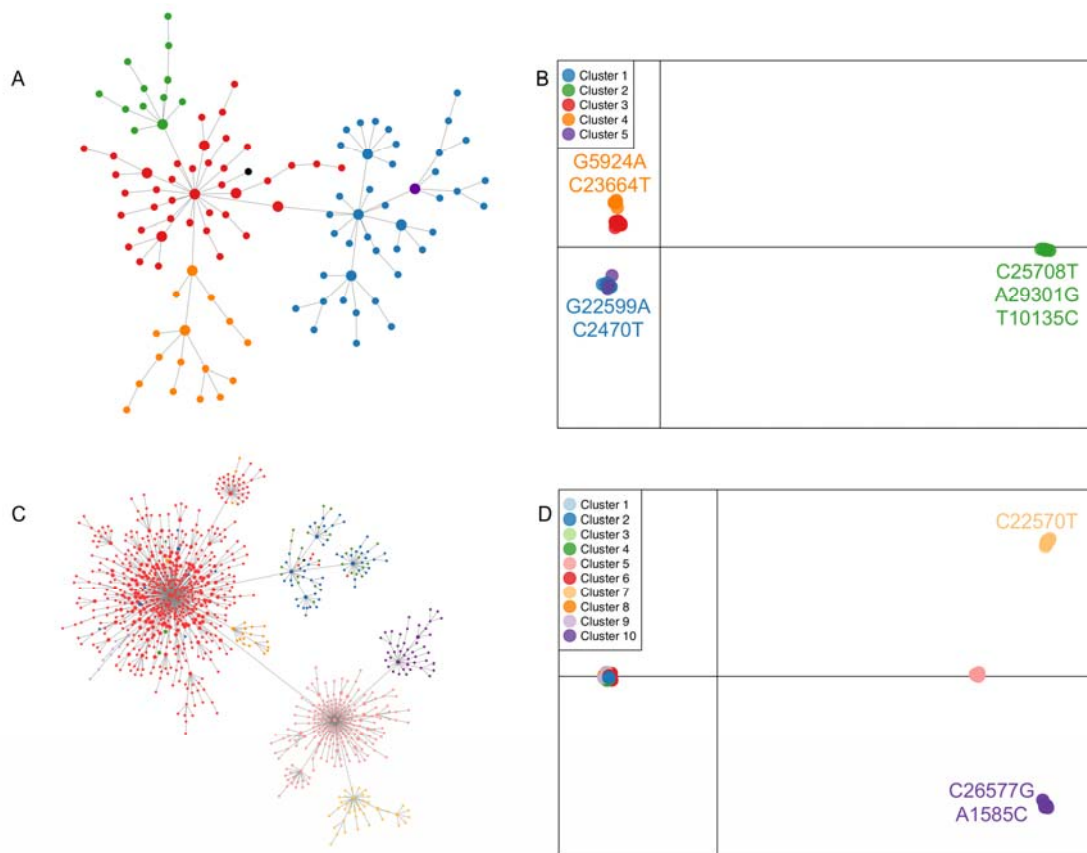
163

164 **REFERENCES**

- 165 1. O'Toole Á, Scher E, Underwood A, Jackson B, Hill V, McCrone JT, et al.
166 Assignment of epidemiological lineages in an emerging pandemic using the pangolin
167 tool. *Virus Evol.* 2021;7:veab064.
- 168 2. Aksamentov I, Roemer C, Hodcroft EB, Neher RA. Nextclade: clade assignment,
169 mutation calling and quality control for viral genomes. *J Open Source Softw.*
170 2021;6:3773.
- 171 3. Ling Y, Cao R, Qian J, Li J, Zhou H, Yuan L, et al. An interactive viral genome
172 evolution network analysis system enabling rapid large-scale molecular tracing of
173 SARS-CoV-2. *Sci Bull (Beijing).* 2022;67:665–9.
- 174 4. Ideker T. Cytoscape: a software environment for integrated models of biomolecular
175 interaction networks. *Genome Res.* 2003;13:2498–504.

- 176 5. Jombart T. adegenet: a R package for the multivariate analysis of genetic markers.
177 Bioinformatics. 2008;24:1403–5.
- 178 6. R Core Team. R: A language and environment for statistical computing. R
179 Foundation for Statistical Computing, Vienna, Austria. 2019. URL [https://www.R-](https://www.R-project.org/)
180 [project.org/](https://www.R-project.org/).
- 181 7. Kumar S, Tamura K, Nei M. MEGA: molecular evolutionary genetics analysis
182 software for microcomputers. Bioinformatics. 1994;10:189–91.
- 183 8. Yang HC, Wang JH, Yang CT, Lin YC, Hsieh HN, Chen PW, Liao HC, Chen CH,
184 Liao JC. Subtyping the major SARS-CoV-2 variants reveals different transmission
185 dynamics. BioRxiv 2022. 10.1101/2022.04.10.486823
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200

201 **FIGURES**



202

203

204 Figure 1: Genetic subdivision observed in Omicron lineage of SARS-CoV-2 based on

205 haplotype network A) BA.1 lineage and C) BA.2 lineage. Panels B and D depict observed

206 genetic subdivision based on DAPC analysis for BA.1 and BA.2 lineages respectively. The

207 Wuhan sequence is denoted in black colored node in panels A and C. Private mutation if any

208 are depicted in the DAPC plot. A private mutation must be present in at least 70% of

209 sequences within the cluster.

210

211

212

213

214 **SUPPLEMENTARY MATERIAL**

215 **TABLES**

216 Table S1: Details of the SARS-CoV-2 genome sequences obtained from GISAID.

217

218 Table S2: Tajima's D estimate for the various clusters identified using DAPC. Values in bold
219 indicate significant demographic expansion or selection.

220

221 **FIGURES**

222

223 Figure S1: Distribution of the number of mutations for each open reading frame of SARS-
224 CoV-2 sequences belonging to omicron lineage analyzed in this study. We characterized the
225 mutation in reference to the Wuhan genome.

226

227 Figure S2: Frequency distribution of the number of A) transitions, and B) transversion for
228 each open reading frame of SARS-CoV-2 sequences belonging to omicron lineage analyzed
229 in this study. We characterized the mutation in reference to the Wuhan genome.

230