

1        **Dynamic tracking of variant frequencies depicts the evolution of mutation**  
2                                **sites amongst SARS-CoV-2 genomes from India**

3

4                                Shenu Hudson B.<sup>1</sup>, Vaishnavi Kolte<sup>1</sup>, Azra Khan, Gaurav Sharma\*

5

6        <sup>1</sup> Equal author

7        \* Corresponding author

8

9

10    **Affiliation:** Institute of Bioinformatics and Applied Biotechnology (IBAB), Bengaluru,  
11    Karnataka, India-560100

12

13    E-mail: [gauravsharma@ibab.ac.in](mailto:gauravsharma@ibab.ac.in)

14

15    **Keywords:** COVID-19; SARS; Coronavirus; Genomics; Bioinformatics; Epidemiology

16

17

18

19 **Abstract:**

20           With the exponential spread of the COVID-19 pandemic across the world within the  
21 twelve months, SARS-CoV-2 strains are continuously trying to adapt themselves in the host  
22 environment by random mutations. While doing so, some variants with evolutionary advantages  
23 such as better human to human transmissibility potential might get naturally selected. This short  
24 communication demonstrates how the mutation frequency patterns are evolving in 2,457 SAR-  
25 CoV-2 strains isolated from COVID-19 patients across diverse Indian states. We have identified  
26 19 such variants showing contrasting mutational probabilities in the span of seven months. Out  
27 of these, 14 variants are showing increasing mutational probabilities suggesting their propagation  
28 with time due to their unexplored evolutionary advantages. Whereas mutational probabilities of  
29 five variants have significantly decreased in June onwards as compared to March/April,  
30 suggesting their termination with time. Further in-depth investigation of these identified variants  
31 will provide valuable knowledge about the evolution, infection strategies, transmission rates, and  
32 epidemiology of SARS-CoV-2.

33

34 **Introduction:**

35           Since the emergence of SARS-CoV-2 in Wuhan, China in December 2019, the infection  
36 has spread at a menacing rate throughout the world. As of today, over 70 million active cases of  
37 COVID-19 and 1.5 million deaths are telling the horror story of this debilitating virus. Since the  
38 whole genome sequencing of the SARS-CoV-2 Wuhan Hu-1 strain earlier this year, more than  
39 250,000 genome sequences<sup>1-3</sup> have been deposited in open-source platforms such as GISAID,  
40 NCBI Virus, etc. Consortiums such as Nextstrain are providing valuable, unprecedented insights  
41 into the demography and epidemiology of SARS-CoV-2 strains along with information to  
42 supervise drug/vaccine design and discovery.

43           Coronaviruses encode a 3'-5' exoribonuclease (Nsp14) that proofreads RNA<sup>4</sup>, which tries  
44 to maintain genome fidelity and control variations, although even with this, mutation rates  
45 amongst these viruses are very high. However, during the past twelve months of the ongoing  
46 COVID-19 pandemic, worldwide human-to-human transmission has enabled SARS-CoV-2 to  
47 accumulate numerous genetic variations<sup>5</sup>. Several mutations among these are non-significant  
48 while other mutations could be beneficial for the survival of the virus, its infecting capability,  
49 and transmission<sup>6-9</sup>. Given that beneficial mutations could be naturally selected in wider  
50 populations, studying SARS-CoV-2 genomic variants and their tracking with time might help us  
51 in understanding viral evolution, behavior, and infection trajectory.

52           In this study, we have tracked and identified several mutational sites for month-wise  
53 (March to September; seven months) separated SARS-CoV-2 genome datasets isolated from  
54 Indian COVID-19 patients to distinguish which variants are getting naturally selected to  
55 propagate further and the ones which are being dismissed from the population.

56

57

58 **Methods:**

59 2,457 complete genomes of SARS-CoV-2 isolated from Indian patients till September 30,  
60 2020, were extracted from GISAID on November 1, 2020 (**Table S1**). From September 30 to  
61 date, only eight Indian SARS-CoV-2 genome sequences have been submitted to GISAID, owing  
62 to this, we have not used genome data after September month in this study. Based on the sample  
63 collection date of each strain, we distributed them into seven categories: 'March' (this included  
64 the strains collected in January, February, and March), April, May, June, July, August, and  
65 September. The genome of the hCoV-19/Wuhan/WH01/2019 strain (EPI\_ISL\_406798) was used  
66 as a reference throughout this analysis<sup>3</sup>. For each category, we generated multiple sequence  
67 alignment and identified the percentage of all four nucleotides (along with gap (-) and other non-  
68 standard nucleotides) at each site. Sites having >2% gaps or non-standard nucleotides were not  
69 considered in the analysis. The frequency of each nucleotide in the alignment was calculated and  
70 a ratio was determined with the frequency of the nucleotide in the reference genome. Mutation  
71 frequency/probability was defined as the ratio of the frequency of the nucleotide at any site to the  
72 frequency of the nucleotide present in the reference sequence at the same aligned site. A  
73 hierarchical clustering based heatmap of each nucleotide loci was generated using mutational  
74 probabilities within each category using the hclust function in R. Simultaneously, trend plots  
75 were also generated for all identified clusters using ggplot.

76

## 77 **Results:**

78 All 2,457 SARS-CoV-2 strains isolated from the diverse landscape of India during  
79 March-September 2020 were categorized into month categories (**Table S1**) and aligned with the  
80 reference WH01 genome. This study recognized 268 sites with mutation probability ranging  
81 from 2 to 97.99% in at least one month. We identified that in most of these sites, there were  
82 negligible variations (>2% and <4%) amongst different month categories. Therefore, to identify  
83 the critical variations in our downstream analysis, we used a 4% minimum mutational probability  
84 score in at least one-month category and found 118 sites encompassing this criterion.  
85 Accordingly, we identified 36, 33, 32, 36, 37, 37, and 50 highly mutating sites amongst the  
86 March, April, May, June, July, August, September categories, respectively. Of these, 11 sites  
87 were showing significant mutation probabilities in all seven months, 3 in six of the months, 8 in  
88 four of the months, 11 in three categories, 16 in two, and 69 were unique in just a single month  
89 category. Finally, the mutational probabilities of these 118 sites across all time points (month-  
90 wise) were visualized using a clustered heatmap where six clusters were obtained with 4, 3, 2, 5,  
91 5, and 99 sites per cluster (**Figure 1AB, Table S2**). The largest cluster with 99 sites did not show  
92 an upward or downward longitudinal trend, therefore, it might be classified as a “neutral” cluster  
93 (**Figure 1AB**), whereas 19 sites distributed in other five clusters showed significant variations,  
94 which are discussed below.

95 **Cluster one** included four loci where the mutational probabilities were lowest in March  
96 (2-35%) with a radical increase to >90% in June and onwards, indicative of positive selection in  
97 the population (**Figure 1AB, Table S2**). The first site, C3037T is in the nsp3 protein-coding  
98 region of the orf1a gene, C14408T is in RdRp (RNA-dependent RNA polymerase) protein-  
99 coding region within the orf1ab gene, the third site (A23403G) is in the spike protein-coding  
100 gene, and the fourth site (C241T) is in the noncoding 5' UTR region. C3037T is a synonymous  
101 mutation (F106F in nsp3 protein) and does not lead to any major changes in protein structure and  
102 function. On the other hand, we noticed that C14408T mutation causes a nonsynonymous change  
103 (corresponding to P125L in RdRp) which has also been observed in genomes isolated from  
104 different continents suggested to change the rigidify of RdRp structure<sup>10</sup>. The third site,  
105 A23403G, is well known to cause a D614G mutation in the spike protein, which interferes in  
106 domain S1-S2 interaction. This has been suggested to cause substantial conformational shifts in

107 spike protein that in turn enhance SARS-CoV-2 infectivity while retaining sensitivity to  
108 antibodies that target the receptor-binding domain <sup>6,7</sup>. Sites identified in the **second cluster**  
109 started showing mutations in April and in each month their frequency is expanding, reaching up  
110 to 60-80% of the total population in August/September. This cluster consists of three sites, of  
111 which one is synonymous mutations (**Table S2**) while the other two are non-synonymous and  
112 worth exploring further. All three sites G28881A, G28882A, and G28883C are located  
113 consecutively in the nucleocapsid gene and lead to changes in amino acids R203K, R203R, and  
114 G204R, respectively.

115 **Cluster three** shows a trend like cluster two however the increase in percentage  
116 representation is slow and goes up to 35-60% as compared to 60-80% in the latter. It consists of  
117 two sites only i.e., a synonymous mutation in C313T [nsp1 in ORF1a; L16L] and a non-  
118 synonymous mutation in C5700A [nsp3 in ORF1a; A994D]. The **fourth cluster** consists of five  
119 sites that had high mutational probabilities during early infection months, *i.e.*, March and April,  
120 and decreased mutational probabilities in May and June, and further decreasing to zero in August  
121 and September, indicating that these variants are not being selected for further propagation. Two  
122 out of these sites, C13730T [which results in an A4489V change in ORF1ab (A97V in RdRp  
123 protein)] and C23929T [which results in a Y789Y change in spike protein], either maintain the  
124 general chemical nature of the amino acid (Alanine to Valine change) or are fully synonymous.  
125 However, the other three sites, C6312A, G11083T, and C28311T, result in non-synonymous  
126 changes in the nsp3 region (T1198K) of the ORF1a (T2016K), the nsp6 region (L37F) of ORF1a  
127 (L3606F), and the nucleocapsid (P13L) protein, respectively. Based on trends, we can  
128 hypothesize that the later three non-synonymous mutations were not selected by the virus  
129 population for further propagation on account of their putative lesser efficiency in infection or  
130 other types of fitness disadvantages.

131 The **fifth cluster** consists of five sites where the mutational probabilities were making a  
132 bell-shaped curve across all month-wise datasets meaning that in March and April along with  
133 August and September, the mutational frequencies were in the range of 2-10%, however, we  
134 could see >30% mutational representation in June and July months. G25563T site is a part of the  
135 orf3a gene and leads to Q57H non-synonymous mutation. Similarly, another non-synonymous  
136 mutation was identified in C28854T that leads to S194L variation in Nucleocapsid protein.

137 Other loci include synonymous mutations at C18877T [nsp14], C22444T [S], and C26735T [M],  
138 which cause insignificant variations as indicated by L280L, D294D, and Y71Y, respectively.  
139 Most of the **sixth cluster** sites have small variations ranging from 0-10% mutational  
140 representations in one of the months. However, in a few sporadic sites, the representation was  
141 10-40% in only one of the months and in other categories, the representation was non-significant  
142 suggesting the absence of any pattern. Finally, considering the ongoing trends of these variants,  
143 we hypothesize that sites identified in clusters 1, 2, and 3 must be getting selected for  
144 propagation owing to some unknown fitness advantages. Similarly, as the representation of sites  
145 identified in clusters 4 and 5 was higher earlier and at its lowest in the end months, we argue that  
146 these sites are not getting preferred during viral population selection.

147 To counter evolutionary pressure, viruses, akin to other living beings, continuously  
148 mutate their genetic material to improve their infection strategies, resistance potential to antiviral  
149 therapies, and transmission rate <sup>11-13</sup>. Most of these random mutations are synonymous or  
150 functionally insignificant. However, a few non-synonymous mutations might give an extra  
151 advantage to the virus in its faster transmission, additional infection severity, or higher resistance  
152 against antiviral vaccines/treatments along with other fitness advantages. Overall, this study  
153 suggests that all identified mutations are not evenly distributed across the virus population during  
154 different timeframes; however, some loci are more prone to propagate and some get terminated  
155 with time. As several variants in clusters 1, 2, and 3 have higher mutational probability in  
156 August/September as compared to March/April, understanding the consequences of those  
157 propagating variants in terms of infection and epidemiology will be of great importance.  
158 Exploring and recognizing this information might prove helpful for drug and vaccine  
159 development. Some reports have already shown the rapid increase of a non-synonymous  
160 (D614G) variant across the globe that might have facilitated increased human-to-human  
161 transmission <sup>6,7</sup>. In-depth investigations of variants identified in this study will provide newer  
162 insights into the evolution and fitness advantages acquired by SARS-CoV-2.

163

164 **Declarations:**

165

166 **Ethics approval and consent to participate**

167 Not applicable

168

169 **Consent for publication**

170 Not applicable.

171

172 **Availability of data and materials**

173 All data and materials are included in this published article.

174

175 **Competing interests**

176 Authors declare that they have no conflicts of interest.

177

178 **Funding**

179 This work was supported by the DST-INSPIRE (Department of Science and Technology)

180 Faculty Award to GS from the Government of India.

181

182 **Authors' contributions**

183 GS conceived and designed the study and wrote the paper. SHB, VK, AK, and GS

184 performed the study and analyzed the data. All authors read and approved the final manuscript.

185

186 **Acknowledgments**

187 GS would like to thank the Department of Science and Technology (DST) and IBAB,

188 Bengaluru for financial and infrastructure support. The authors would like to thank his colleague

189 Dr. Shruthi S. Vembar for constructive criticism of the manuscript. The authors would also like

190 to thank Dr. R. Srivatsan for his help in statistical analysis. We would also like to thank the

191 GISAID database for allowing us to access the genome sequences for this scientific research.

192



193 **Legends for Figures:**

194 **Figure 1:** Dynamic tracking of mutational frequencies within genomes extracted from Indian  
195 samples in March, April, May, June, July, August, and September. **A) Heatmap of mutational**  
196 **probabilities. B) Trend plots of average mutational probabilities within each cluster.** For  
197 each cluster, n represents the number of sites. Time period on X-axis denotes month categories.

198

199 **Legends for Supplementary Tables:**

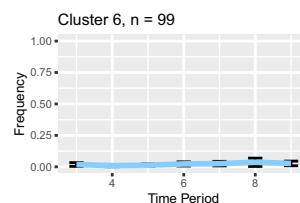
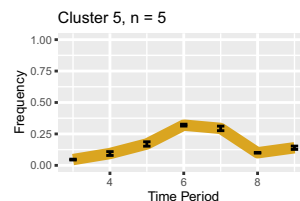
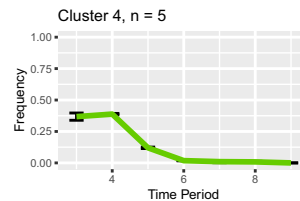
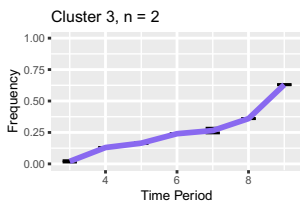
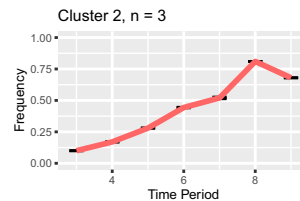
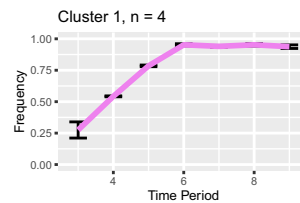
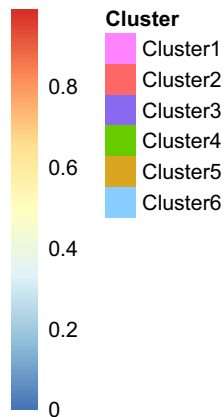
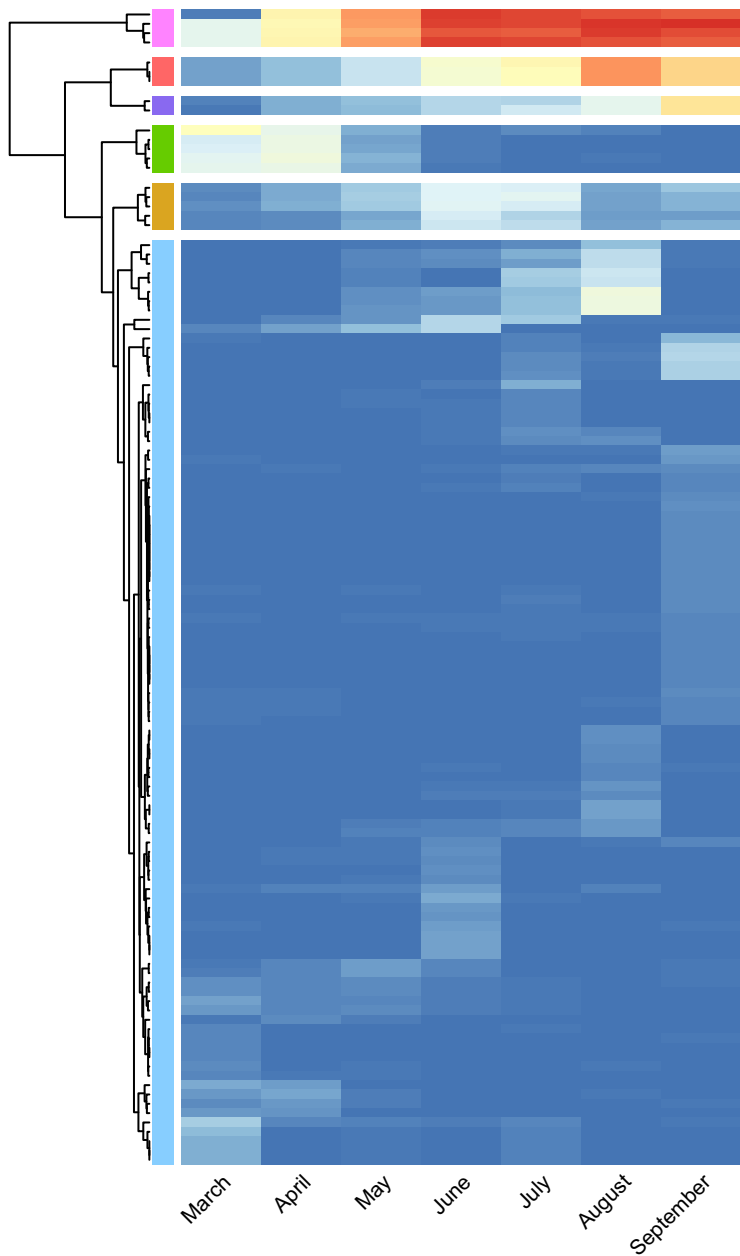
200 **Table S1:** Location-based distribution of 2,457 Indian SARS-CoV-2 strains analyzed in this  
201 study.

202 **Table S2:** Mutation probability scores and other relevant metadata for 268 identified mutation  
203 sites at Genome, Codon, and protein level. Heatmap-based cluster information using 4%  
204 minimum mutational probability data (118 sites) is also provided in the table. NA represents  
205 “Not Applicable” for cluster information as those sites were not considered for heat map-based  
206 clustering.

207

208 **References:**

- 209 1. Lescure FX, Bouadma L, Nguyen D, et al. Clinical and virological data of the first cases of COVID-19  
210 in Europe: a case series. *Lancet Infect Dis.* 2020;0(0). doi:10.1016/S1473-3099(20)30200-0
- 211 2. Zhu N, Zhang D, Wang W, et al. A novel coronavirus from patients with pneumonia in China, 2019. *N*  
212 *Engl J Med.* 2020;382(8):727-733. doi:10.1056/NEJMoa2001017
- 213 3. Fuk-Woo Chan J, Kok K-H, Zhu Z, et al. Genomic characterization of the 2019 novel human-  
214 pathogenic coronavirus isolated from a patient with atypical pneumonia after visiting Wuhan. 2020.  
215 doi:10.1080/22221751.2020.1719902
- 216 4. Sevajol M, Subissi L, Decroly E, Canard B, Imbert I. Insights into RNA synthesis, capping, and  
217 proofreading mechanisms of SARS-coronavirus. *Virus Res.* 2014;194:90-99.  
218 doi:10.1016/j.virusres.2014.10.008
- 219 5. Koyama T, Weeraratne D, Snowdon JL, Parida L. Emergence of Drift Variants That May Affect  
220 COVID-19 Vaccine Development and Antibody Treatment. *Pathogens.* 2020;9(5):324.  
221 doi:10.3390/pathogens9050324
- 222 6. Yurkovetskiy L, Pascal KE, Tomkins-Tinch C, et al. SARS-CoV-2 Spike protein variant D614G  
223 increases infectivity and retains sensitivity to antibodies that target the receptor binding domain.  
224 *bioRxiv.* July 2020:2020.07.04.187757. doi:10.1101/2020.07.04.187757
- 225 7. Korber B, Fischer WM, Gnanakaran S, et al. Tracking changes in SARS-CoV-2 Spike: evidence that  
226 D614G increases infectivity of the COVID-19 virus. *Cell.* July 2020. doi:10.1016/j.cell.2020.06.043
- 227 8. Ryder SP. Analysis of Rapidly Emerging Variants in Structured Regions of the SARS-CoV-2  
228 Genome. *bioRxiv.* June 2020:2020.05.27.120105. doi:10.1101/2020.05.27.120105
- 229 9. Simmonds P. Rampant C→U Hypermutation in the Genomes of SARS-CoV-2 and Other  
230 Coronaviruses: Causes and Consequences for Their Short- and Long-Term Evolutionary Trajectories.  
231 *mSphere.* 2020;5(3). doi:10.1128/msphere.00408-20
- 232 10. Eskier D, Karakulah G, Suner A, Oktay Y. RdRp mutations are associated with SARS-CoV-2  
233 genome evolution. *bioRxiv.* May 2020:2020.05.20.104885. doi:10.1101/2020.05.20.104885
- 234 11. Peck KM, Lauring AS. Complexities of Viral Mutation Rates. *J Virol.* 2018;92(14).  
235 doi:10.1128/jvi.01031-17
- 236 12. Gorbalenya AE, Enjuanes L, Ziebuhr J, Snijder EJ. Nidovirales: Evolving the largest RNA virus  
237 genome. *Virus Res.* 2006;117(1):17-37. doi:10.1016/j.virusres.2006.01.017
- 238 13. Lynch M, Ackerman MS, Gout JF, et al. Genetic drift, selection and the evolution of the mutation  
239 rate. *Nat Rev Genet.* 2016;17(11):704-714. doi:10.1038/nrg.2016.104



A.

B.