1    **ORIGINAL ARTICLE**

2    **Phylodynamic Pattern of Genetic Clusters, Paradigm Shift on Spatio-temporal Distribution of**

3    **Clades, and Impact of Spike Glycoprotein Mutations of SARS-CoV-2 Isolates from India**

4

5    Srinivasan Sivasubramanian[1], Vidya Gopalan[1], Kiruba Ramesh[1], Padmapriya Padmanabhan[1],

6    Kiruthiga Mone[1], Karthikeyan Govindan[1], Selvakumar Velladurai[1], Kaveri Krishnasamy[1], Satish

7    Srinivas Kitambi[2,*]

8

9    [1] State Viral Research and Diagnostic Laboratory (VRDL) Unit, Department of       Virology,   King

10   Institute of Preventive Medicine and Research, Chennai, Tamil Nadu,       India – 600 032.

11   2:Institutet for Healthcare Education and Translational Sciences (IHETS), Hyderabad, Telengana

12   -500026.

13

14

15   **\*Corresponding author.       Dr. Satish Srinivas Kitambi**

16   Institute for Healthcare Education and Translational Sciences (IHETS)

17   10-2-311, Plot 187, Str4, Cama Manor, West Marredpally, Telengana 500026.

18   **Email address: satish.kitambi@klife.info**

19

20   **Running Title:** Analyses of SARS-CoV-2 isolates from India

21   **Keywords**: SARS-CoV-2, COVID-19, Spike protein, Clade, Phylogeny, Mutations, India

22

23    **Abstract**

24    **Background:** The COVID-19 pandemic is associated with high morbidity and mortality, with

25    the emergence of numerous variants. The dynamics of SARS-CoV-2 with respect to clade

26    distribution is uneven, unpredictable and fast changing. **Aims:** Retrieving the complete genomes

27    of SARS-CoV-2 from India and subjecting them to analysis on phylogenetic clade diversity,

28    Spike (S) protein mutations and their functional consequences such as immune escape features

29    and impact on infectivity. **Methods:** Whole genome of SARS-CoV-2 isolates (n=4,326)

30    deposited from India during the period from January 2020 to December 2020 is retrieved from

31    GISAID and various analyses performed using *in silico* tools. **Results:** Notable clade dynamicity

32    is observed indicating the emergence of diverse SARS-CoV-2 variants across the country. GR

33    clade is predominant over the other clades and the distribution pattern of clades is uneven.

34    D614G is the commonest and predominant mutation found among the S-protein followed by

35    L54F. Mutation score prediction analyses reveal that there are several mutations in S-protein

36    including the RBD and NTD regions that can influence the virulence of virus. Besides, mutations

37    having immune escape features as well as impacting the immunogenicity and virulence through

38    changes in the glycosylation patterns are identified. **Conclusions:** The study has revealed

39    emergence of variants with shifting of clade dynamics within a year in India. It is shown uneven

40    distribution of clades across the nation requiring timely deposition of SARS-CoV-2 sequences.

41    Functional evaluation of mutations in S-protein reveals their significance in virulence, immune

42    escape features and disease severity besides impacting therapeutics and prophylaxis.

43

44

45

## INTRODUCTION

Analyses of global and Indian SARS-CoV-2 genome sequences (as on December 2020) have revealed that the virus has differentially distributed into at least 10 clades and is continuously evolving.[1] The S-protein of SARS-CoV-2 targets angiotensin-converting enzyme 2 (ACE2) receptor for its entry into target cells. This protein is the major focus of the vaccine development platforms. Changes in the O- and N- linked glycosylation patterns of the S protein have an impact on the immunogenicity and virulence of the virus. Hence, it is important to closely monitor antigenic evolution of the S-protein in the circulating viruses. In this study, we retrieved complete genomes of SARS-CoV-2 from India during the whole year period from GISAID and subjected them to the studies on clade analyses and clade distribution pattern covering all states of the country. Further, mutations in various regions of S-protein, mutation frequency, glycosylation patterns, and the effects on protein structure, immunity and virulence were analysed.

## METHODS

### Genome data retrieval, phylogenetic and clade analysis

A total of 4326 annotated SARS CoV-2 whole genome sequences (WGS) from various parts of India deposited as on 31st December 2020 in Global Initiative on Sharing All Influenza Data (GISAID)(https://www.gisaid.org/) were retrieved. Sequences were aligned using MAFFT (Multiple Alignment using Fast Fourier Transform) with SARS-CoV-2 Wuhan-Hu-1 strain (NC_045512.2) and GISAID reference sequence (EPI_ISL_402124)[2] used as reference. The Nextclade-Nextstrain pipeline (https://clades.nextstrain.org/) was used for studies on phylogenetic analysis and clustering patterns of the S gene.[3] Further, the Average evolutionary

69  divergence was estimated using Kimura-2 parameter model. Evolutionary analyses and

70  phylogenetic tree construction were performed using MEGA-X.[4]

71

72  **Frequency and functional evaluation of variants**

73  Frequencies and amino-acid variants were analyzed using COVID CG and Tracking mutation

74  tools (source GISAID) respectively. Functional consequences were predicted using tools like

75  SIFT (Sorting Intolerant from Tolerant)

76  (https://sift.bii.astar.edu.sg/www/SIFT_seq_submit2.html),[5] PROVEAN (Protein Variation

77  Effect Analyzer) (http://provean.jcvi.org/seq_submit.php)[5] and PolyPhen-2 (Polymorphism

78  Phenotyping v2) (http://genetics.bwh.harvard.edu/pph2).[6] A SIFT score of 0.0 to 0.05 indicates

79  a deleterious effect. The functional effects of protein variants were assessed using the

80  PROVEAN web server, using a default threshold value of −2.5 and the values below and above

81  the threshold value are considered as deleterious and tolerant respectively. A PolyPhen threshold

82  scores of > 0.908, >0.446 and ≤ 0.908 and ≤ 0.446 are interpreted as "Probably Damaging",

83  "Possibly Damaging" and "Benign" respectively. ESC_Comprehensive resource of immune

84  escape variants in SARS-COV-2 was used to detect the escape mutants in S-protein

85  (http://clingen.igib.res.in/esc/).[7]

86

87  **RESULTS**

88  The retrieved WGS were found to be classified under 7 clusters according to GISAID Clade

89  identification [Figure 1].[8] It was observed that the predominant cluster encompassed 1,755

90  (40.56%) of genomes that fell under the GR clade [Figure 1a]. Though this clade was

91  represented by samples derived from various states across India, Maharashtra (n = 922) and

92   Telangana (n = 492) states had the maximum numbers followed by Karnataka (n = 102) [Figure

93   1b]. The major clade GR was followed by clades G (942; 21.77%), O (783; 18.1%), GH (737;

94   17.03%), S (82; 1.9%), L (25; 0.58%) and V (3; 0.07%). The clade G was mainly represented by

95   samples from Maharashtra (n = 277), Gujarat (n = 215) and West Bengal (n = 152). The O clade

96   is prevalent in all states of the country. Gujarat state accounts for the highest number of samples

97   under GH clade [Figure 1b]. States such as Andhra Pradesh, Punjab and Rajasthan submitted a

98   smaller number of sequences and the clade diversity pattern could not be clearly deciphered.

99

100   The viruses belonging to L, S and O clades were prevalent during the initial months (January

101   to February, 2020) [Figure 1c]. During the starting of the pandemic (March to April), O clade

102   was predominant followed by G, GR, GH, S and L. It is noteworthy that the distribution of S and

103   L clades were drastically reduced during this period and the strains belonged to clades O, S, L

104   and V were remarkably low in numbers during the progress of pandemic. From May to October,

105   GR clade is predominant but becomes second to GH clade during November and December.

106   Notably, the O clade was slowly dominated by GR, G and GH clades in different states during

107   the course of pandemic and there was almost near to complete absence of O clade during

108   November and December. The phylogenetic tree depicting clade diversity throughout the year

109   shows that GR is the dominant clade over the others [Figure 2]. These results suggested

110   spatiotemporal clade diversity and a paradigm shift in phylodynamics of clade distribution.

111

112   **Mutation analysis of spike protein from Indian strains**

113   A total of 557 amino acid substitution mutations were found in S-protein among the 4,326 Indian

114   strains [Supplementary File 1]. There were 333 and 215 mutations present in the S1 and S2

115   domains respectively with the highest number of mutations in the N-terminal domain (NTD; 211

116   mutations) followed by the RBD (63 mutations) [Table 1]. Nine mutations are identified in

117   signal peptide, which is not the component of mature S protein. Among these 557 mutations,

118   D614G was present in 79.99% (n = 3461) of Indian strains followed by L54F (n=111, 2.57%)

119   isolates. The other prominent mutation sites were: Q677 (72), P681 (54), P812 (40), A771 (34),

120   Q675 (30), and L5 (26) [Figure 3]. Besides, 11 types of mutations are found in the 8 sites of

121   highly conserved protease cleavage region (from 675 to 692 of S1 and S2 domains) of the

122   protein. L18F, H69del, V70del, D138Y and Y144del mutations were observed in NTD of S-

123   protein of few isolates and these mutations could enhance the surface electropositivity of the S-

124   protein and thereby facilitating the adhesion of virus to negatively charged lipid raft gangliosides

125   of host cells.[9] It is also observed that two of the study variants possess H69del, V70del and

126   Y144del in NTD and N501Y in RBD suggesting the improved affinity as well as adhesive

127   properties of S-protein due to the concomitant mutations in both regions that synergistically

128   promote virus host interaction.

129

130        The frequency of amino acid mutations in S-protein was analyzed using COVID-19 CoV

131   Genetics browser (source: GISAID), and the results showed that non-synonymous mutations

132   were scattered across the S-gene with region specific varying frequency [Supplementary File 2].

133   Figure 4 shows prevalent mutations such as D614G, Q677H and P681H originated during

134   March, April and July respectively and their appearance was observed till the end of the year

135   2020. On contrary, L54F as well as K77M and P812L mutations emerged during April and June

136   respectively but absent after few months of their appearance.

137

138        Many amino acid mutations were observed to be region specific namely F32Y, T33K and

139    G35Q mutations (in Karnataka); T29I and P681H (Maharashtra); and L7S, L54F, R78M,

140    Q690H, A701T and A879S (Gujarat). These mutations were absent from other states indicating

141    that these mutations might not spread to other states possibly due to effective implementation of

142    lockdown measures throughout the country. Some distinct amino acid variants were observed in

143    Gujarat and Maharashtra (G181A) and V622F in Telangana and Orissa. There were 12

144    premature stop codons and 8 deletion mutations present in different positions of various S-gene

145    sequences. More than one mutation type can be observed at the same position in the protein. For

146    instance, amino acid A to V, E, S, or K, at position 27, A to G, T, S, or V at position 222. Among

147    the total 419 mutation sites in S-protein of Indian isolates, 114 sites carry more than one

148    mutation. It is noteworthy that there were 190 distinct mutation events that occurred in India first

149    time; among them, 115 mutation events were confined only to India and the rest of 75 mutations

150    were subsequently identified in various countries or occurred independently at different

151    geographical regions across the world [Supplementary File 2].

152

153    **Immune escape mutations in spike protein**

154    The analyses showed 11 and 17 immune escape mutations in the NTD and RBD of S-protein

155    respectively [Supplementary File 3]. L18F, T19A, D80N, D138Y, Y144del, Y145del, K147E,

156    N148S, W152L, Q218H and S255F were found in NTD, and among them, L18F, Y144del,

157    Y145del and N148S and W152L were shown to display resistance to neutralizing antibodies.

158    Among the mutations in RBD, R346K, N440K, G446V, N450K, V483F, E484K, E484Q, F490S

159    and S494P also showed change in ACE2 binding to the extent of 75% to 90%.[10,11] Variants

160    identified with mutations at sites such as E484 (E484Q), F490 (F490S), Q493 (Q493STOP), and

161    S494 (S494P) in the RBD are presumed to have immune escape features.[12]

162

163        Mutations in regions of S-protein other than RBD also can show resistance to antibody

164    and are identified in the present study. It is noteworthy that single amino acid changes such as

165    Y145del, F490S, A831S and double amino acid changes including D614G+A879S,

166    D614G+A879T, and D614G+M1237I were reported to be resistant to convalescent sera or these

167    mutations could confer the S protein monoclonal antibody resistance, whereas V367F of the

168    RBD was reported to have increased sensitivity to neutralizing antibodies.[13] Other mutations

169    M153I, S254F, and S255F identified in the study are found to reduce the affinity between S-

170    protein and antibodies.[14]

171

**Mutations affecting glycosylation patterns**

172

173    Analysis of both N- linked (NGS) and O-linked glycosylation sites (OGS) was performed for S-

174    protein of 4326 isolates. Among the total 22 and 26 amino acid sites of S-protein carrying with

175    NGS and OGS moieties respectively, it was observed that 7 and 9 of these sites were found to

176    possess mutations that resulted in loss of glycosylation moiety [Figure 5]. All except one variants

177    possessed only one amino acid glycosylation site change either NGS or OGS. One variant

178    (EPI_ISL_479737) had lost both OGS and NGS sites due to mutations such as T602L and

179    N603Y [Table 2]. There were two NGS present in RBD without any mutation; among the four

180    OGS in RBD, only one glycosylation mutation (T323I) was observed.

181

182

**Functional evaluation of the S protein mutations**

Among the total 557 amino acid mutations of S-protein, 531 mutations were taken up for score prediction studies and the remaining 26 mutations observed either as stop codons (STOP) or deletions (del). SIFT score predicted 124 mutations to be deleterious and other mutations to be neutral. Also, PROVEAN score predicted 63 mutations to be deleterious whereas POLYPHEN-2 predicted 213 mutations that could display probably damaging effect. Only 41 amino acid mutations were predicted to result in potentially deleterious functional consequences by all three of the mutation score prediction tools [Supplementary File 4].

**Phylogenetic analysis of spike protein**

Only 250 S-protein sequences constituting unique mutations were selected for phylogenetic tree construction, and the analysis showed that there was high degree of heterogeneity with multiple clusters and sequences were highly diverged from the reference sequence [Figure 6]. Estimates of Average Evolutionary Divergence of sequence pairs comprising 4312 S-genes showed the evolutionary rate as 5.4 x $10^{-4}$ substitution/site/year (s/s/y).

**DISCUSSION**

Continuous monitoring of the virus locally and globally is needed for devising effective measures to handle the pandemic crisis. In this study, we report the molecular epidemiological features of SARS-CoV-2 based on WGS in GISAID deposited from India including the dynamics of clade distribution and diversity, amino acid mutations in S-protein and their impact on virulence, immune evading responses and glycosylation patterns.

206    The study showed that the GR was predominant and was followed by clades G, O, GH, S,

207    L and V. Though there were only four SARS-CoV-2 clades such as 'L', 'S', 'G', and 'V' during

208    the early pandemic phase, swift genetic diversity of the virus and its rapid pace of evolution

209    facilitated GISAID to continuously update the classification by inclusion of three more clades

210    such as GH', 'GR' and 'GV'. Besides, all unclassified sequences of SARS-CoV-2 strains are

211    grouped as 'O. It is observed that there are only few studies on phylogenetic analyses of SARS-

212    CoV-2 from India. A recent study from India reported that the major cluster of SARS-CoV-2

213    was A2a (PANGOLIN lineage B.1/B.1.1/B.1.36) (83%) followed by a distinct A3i clade

214    (PANGOLIN lineage B.6) (11.6%).[5,15] Another phylogenetic study on Indian SARS-CoV-2

215    revealed the presence of four major clades, *i.e.*, 19A (n = 18.4%), 19B (n = 17%), 20A (n =

216    34.43%), 20B (n = 28.3%), and one minor clade 20C (n = 1.9%).[16] These reports suggested that

217    Europe and Southeast Asia as two major routes for introduction of the virus in India followed by

218    local transmission. Both the predominant G and GR are European clades and the strains of these

219    clades possess the D614G mutation on the S-protein which is more infectious.[16,17]

220

221    The month-wise clade distribution analysis showed that L, S and O clades were prevalent

222    in the country during the early phase of pandemic; subsequently, G, GR, and GH clades became

223    prevalent over them. The prevailing clades in the country could be attributed to the early

224    invasion of strains into India through travelers and subsequent mixing of clades. Few reports

225    with minimal sequences deposited till July 2020 only revealed the presence of few clades such as

226    A2a, A3, B and O in India and among them A2a (related to GISAID clade G) was predominant

227    following A3, O and B.[5,17,18] The present study observes ever-changing genetic diversity with

228    intense clade dynamicity of the virus throughout the year.

229

230    Substitution mutations in S protein of all the Indian SARS-CoV-2 sequences were

231    analysed with reference to SARS-CoV-2 Wuhan-Hu-1 strain. The origin of D614G mutation was

232    in China during January, 2020 but the occurrence in India was reported in March and became

233    prevalent afterwards. The first occurrence of L54F was observed in Wuhan in March whereas

234    India reported in April. The protease cleavage region S1/S2 in the S protein is essential for the

235    virus to undergo proteolytic activation of S1 and S2 domains for receptor binding and viral-

236    membrane fusion. The region is highly conserved at sites 685 and 686 where proteolytic

237    cleavage occurs. The study has identified 11 mutations flanking the proteolytic cleavage site.

238    Inferences from the proteolytic cleavage of the S glycoprotein suggest the capability of virus to

239    possess features such as cross species mobility or tropism towards different cells.[19] There are

240    166 mutation sites observed in Asia with 181 mutation types.[20] However, the present study

241    observes that there are 419 mutation sites in the protein with 557 mutation types meaning that

242    several sites in the protein carry more than one mutation type.

243

244    Though D614G is associated with increased infectivity, mutations such as Q239R, T719I,

245    T719S, D839Y, P1263L, mutations in RBD such as I434K and P521S, and D614G+Q675H are

246    reported to have decreased infectivity.[13] Besides, D614G in combination with other mutations

247    such as D614G+L5F (n = 23), D614G+V341I (n = 1), D614G+D936Y (n = 3), D614G+S939F

248    (n = 9) and D614G+S943T (n = 2) in strains of the present study was demonstrated to have

249    increased infectivity compared to Wuhan-1 strain.[13] A recent study has reported that L54F,

250    D614G and V1176F of S-protein, identified in the study, are correlated with severe clinical

251    outcome.[21] It was reported that mutations such as T29I, H49Y, D138Y, E484Q, E484K, A520S,

252   T572I, D614G and H1083Q identified in strains of the study, could increase the stability of S-

253   protein.[6] In contrast, the report suggested that mutations such as L54F, G431S, E471D, G502R,

254   Q506H, P507S, Y508N, E583D and Q675H could weaken the interaction of S-protein with

255   ACE2 receptor; whereas, N440K, E471Q and G504V could improve the binding affinity.

256   Emergence of strains of Variant of Concern (VOC), according to WHO nomenclature, such as

257   Alpha (GISAID clade: GRY), Beta (GH/501Y.V2), Gamma (GR/501Y.V3) and Delta

258   (G/452.V3) as well as Variant of Interest (VOI) such as Eta (G/484K.V3), Iota (GH/253G.V1)

259   and Kappa (G/452R.V3) has been observed during the end of year 2020 and early 2021

260   worldwide. Though few of these highly transmissible variants identified in India late 2020, the

261   sequences of them were submitted to GISAID only in 2021 except two VOC Alpha strains

262   (EPI_ISL_745197 and EPI_ISL_747244). Hence, the study does not report mutations and their

263   features for these variants including the Delta variant that are likely responsible for the

264   substantial surge in cases that began in the Western state of Maharashtra and spread throughout

265   India from Jan, 2021 onwards.[22] This study observed that only 2 WGS of VOC strain (Alpha)

266   from India were available in GISAID in the year 2020 and were taken for analysis.

267

268   Antibodies targeting the RBD are being developed as prophylactics. Determination of

269   mutations in S-protein showing resistance to antibodies is crucial for assessing the antigenic

270   implications of viral evolution. The study has identified immune escape mutations both in NTD

271   and RBD of Indian isolates. Mutations especially in these domains evading the antibody

272   recognition could result in the severity of infection. Presently, most of the SARS-CoV-2 genome

273   is not under positive selection, but if neutralizing antibodies are widely deployed as

274   prophylactics, positive selection pressure that lead to infection-competent viral mutants resulting

275     in resurgence of SARS-CoV-2 infections and pose challenges to prophylactic measures.[11]

276     Virulence of SARS-CoV-2 can be associated with mutations in S-protein such as L18F, H69del,

277     V70del, D138Y and Y144del that confer affinity and adhesive properties for better interaction

278     with host cells through surface electrostatic interaction;[9] besides, these mutations are also

279     reported to evade host immune responses against S-protein.[23] Though the present study

280     particularly focuses on functional features of mutations in S-protein, epistatic interactions

281     involving mutations from other genes can also play a role in clade diversity and spatio-temporal

282     dynamics. Such interactions favor the coevolution of mutations due to selective pressures to form

283     new clades that become dominant. The fitness of mutations in virulence and immune escape

284     features are largely influenced not only by independent mutations in S-protein but also mutations

285     through epistatic interactions. For instance, D614G appears along with 3 other mutations in

286     5'UTR, NSP3 and NSP12 that form G clade.[24] VOC strains forming distinct clades have

287     virulence features contributed by mutations in S gene and other genes.

288

289     Glycosylation of S protein plays a vital role in virulence, S-protein folding, immune

290     sensitivity as well as host immune evasion, and shaping viral tropism.[25] Analysis of both NGS

291     and OGS of the study isolates showed mutations that resulted in loss of glycosylation moiety

292     suggesting the reduced immunogenic potential of S-protein of mutant variants.[26] However,

293     there is no report on the impact of the NGS mutation in the interaction of RBD with ACE

294     receptor. S-protein of SARS-CoV-2 has 22 NGS and several OGS; but, in many strains of this

295     study, several of these sites were lost due to amino acid mutations. There are studies that report

296     absence of mutations at NGS in S-protein. It has been studied that certain mutations incurred in

297     the NGS and OGS increase the stability of the S-protein.[6] Accordingly, in the present study, the

298    observed mutations in the NGS such as N234Y and N603Y and OGS mutations such as S221L,

299    T323I, T602I and T602L are found to stabilize the S-protein. On the contrary, very few

300    mutations at the NGS (N709K) and OGS (T1077I) were found to decrease the stability of S

301    protein.[6] Also, glycosylation mutations such as N149G, N165S, and N709K are reported to

302    increase the sensitivity to neutralizing antibodies and the mutation N234Y is found to reduce the

303    neutralization sensitivity to different set of antibodies. The glycosylation mutation N1074D has

304    been found to decrease the infectivity.[13]

305

306        Functional evaluation of 531 mutations in S-protein from Indian isolates reveals 41

307    amino acid mutations that are predicted to have potential impact on functional consequences. A

308    previous study on Indian SARS-CoV-2 isolates reported scores for D614G mutation in S-protein

309    and several mutations across various proteins with their functional impact.[5] However, the

310    present study reports scores for all mutations occurred in S-protein of Indian isolates that were

311    predicted to be neutral, tolerated, deleterious, benign and probably damaging by means of using

312    mutation score prediction tools. The evolutionary rate of S-gene was estimated to be $5.4 \times 10^{-4}$

313    substitution/site/year (s/s/y) through analysis of 4312 S-genes. Reports suggest that the genome

314    have the evolutionary rate varying in the range between $1.854 \times 10^{-4}$ and $5.63 \times 10^{-3}$s/s/y.[27-29] A

315    study reported that the evolutionary rate for S-gene of SARS-CoV-2 was $1.08 \times 10^{-3}$s/s/y after

316    nine months of pandemic.[30] Another study on Indian SARS-CoV-2 isolates reported the

317    evolutionary rate for S-protein as $3.55 \times 10^{-3}$s/s/y employing sequences of 1376 isolates.[5]

318

319

320

321    **CONCLUSIONS**

322    The study has revealed a rapidly shifting of clade predominance and uneven distribution within a

323    year of the introduction of SARS-CoV-2 in India. The evaluation of S protein reveals the

324    significance of various mutations in virulence, immune escape features and disease severity

325    besides their impact on therapeutics and prophylaxis.

326

327    **Acknowledgements**

332

333

334

335

336

337

338

339

340

341

342

343

344 **Conflicts of interest**

345 None

346

347 **Ethical approval**

348 None.

349

350 **Research Quality and Ethics Statement**

351 The authors of this manuscript declare that this scientific work complies with reporting quality,

352 formatting and reproducibility guidelines set forth by the EQUATOR Network. The authors also

353 attest that this clinical investigation was determined to not require the Institutional Review Board

354 / Ethics Committee review, and the corresponding protocol/approval number is not applicable.

355 We also certify that we have not plagiarized the contents in this submission and have done a

356 Plagiarism Check. We also certify that none of the authors is a member of the Editorial board of

357 the Journal of Global Infectious Diseases.

358

359

360

361

362

363

364

365

366

367    **REFERENCES**

368

369    1.  Sun J, He WT, Wang L, Lai A, Ji X, Zhai X, *et al*. COVID-19: Epidemiology, Evolution,
370        and Cross-disciplinary perspectives. Trends Mol Med 2020;26(5):483-95.

371    2.  Elbe S, Buckland Merrett G. Data, disease and diplomacy: GISAID's innovative
372        contribution to global health. Glob Chall 2017;1(1):33-46.

373    3.  Hadfield J, Megill C, Bell SM, Huddleston J, Potter B, Callender C, *et al*. Nextstrain: real-
374        time tracking of pathogen evolution. Bioinformatics 2018;34(23):4121-3.

375    4.  Kumar S, Stecher G, Li M, Knyaz C, Tamura K. MEGA X: Molecular Evolutionary
376        Genetics Analysis across Computing Platforms. Mol Biol Evol2018;35(6):1547-1549.

377    5.  Banu S, Jolly B, Mukherjee P, Singh P, Khan S, Zaveri L, *et al*. A distinct phylogenetic
378        cluster of Indian SARS-CoV-2 isolates. Open Forum Infect Dis 2020;7(11):ofaa434.

379    6.  Teng S, Sobitan A, Rhoades R, Liu D, Tang Q. Systemic effects of missense mutations on
380        SARS-CoV-2 spike glycoprotein stability and receptor-binding affinity. Brief Bioinform
381        2021;22(2):1239-53

382    7.  Rophina M, Pandhare K, Shamnath A, Imran M, Jolly B, Scaria V. ESC - a
383        comprehensive resource for SARS-CoV-2 immune escape variants. bioRxiv 2021.

384    8.  Alm E, Broberg EK, Connor T,  Hodcroft EB, Komissarov AB, Maurer-Stroh S, *et al*.
385        Geographical and temporal distribution of SARS-CoV-2 clades in the WHO European
386        Region, January to June 2020. Euro Surveill 2020;25(32):2001410.

387    9.  Fantini J, Yahi N, Azzaz F, Chahinian H. Structural dynamics of SARS-CoV-2 variants: A
388        health monitoring strategy for anticipating Covid-19 outbreaks. J Infect 2021;83(2):197-
389        206.

390   10. Van Egeren D, Novokhodko A, Stoddard M, Tran U, Zetter B, Rogers M, *et al*. Risk of
391        rapid evolutionary escape from biomedical interventions targeting SARS-CoV-2 spike
392        protein. PLoS One 2021;16(4):e0250780.

393   11. Greaney AJ, Starr TN, Gilchuk P, Zost SJ, Binshtein E, Loes AN, *et al*. Complete
394        mapping of mutations to the   SARS-CoV-2 spike receptor-binding domain that escape
395        antibody recognition. Cell Host Microbe 2021;29(1):44-57.e9.

396   12. Weisblum Y, Schmidt F, Zhang F, DaSilva J, Poston D, Lorenzi JC, *et al*. Escape from
397        neutralizing antibodies by SARS-CoV-2 spike protein variants. Elife 2020;9:e61312.

398   13. Li Q, Wu J, Nie J, Zhang L, Hao H, Liu S, *et al*. The impact of mutations in SARS-CoV-2
399        Spike on viral infectivity and antigenicity. Cell 2020;182(5):1284-94.e9.

400   14. Chen J, Gao K, Wang R, Wei GW. Prediction and mitigation of mutation threats to
401        COVID-19 vaccines and antibody therapies. Chem Sci. 2021 Apr 13;12(20):6929-6948.

402   15. Jacob JJ, Vasudevan K, Veeraraghavan B, Iyadurai R, Gunasekaran K. Genomic evolution
403        of severe acute respiratory syndrome Coronavirus 2 in India and vaccine impact. Ind J
404        Med Microbiol 2020;38(2):210-2.

405   16. Raghav S, Ghosh A, Turuk J, Kumar S, Jha A, Madhulika S, *et al*.  Analysis of Indian
406        SARS-CoV-2 genomes reveals prevalence of D614G mutation in Spike protein predicting
407        an increase in interaction with TMPRSS2 and virus infectivity. Front Microbiol
408        2020;11:594928.

409   17. Pattabiraman C, Habib F, Rasheed R, Prasad P, Reddy V, Dinesh P, *et al*. Genomic
410        epidemiology reveals multiple introductions and spread of SARS-CoV-2 in the Indian
411        state of Karnataka. PLoS One 2020;15(12):e0243412.

412    18. Biswas NK, Majumder PP. Analysis of RNA sequences of 3636 SARS-CoV-2 collected

413         from 55 countries reveals selective sweep of one virus type. Indian J Med Res

414         2020;151(5):450-8.

415    19. Menachery VD, Dinnon KH, Yount BL, McAnarney ET, Gralinski LE, Hale A, *et al*.

416         Trypsin Treatment Unlocks Barrier for Zoonotic Bat Coronavirus Infection. J Virol

417         2020;94(5):e01774-19.

418    20. Guruprasad L. Human SARS CoV-2 spike protein mutations. Proteins. 2021;89(5):569-

419         76.

420    21. Nagy A, Pongor S, Gyorffy B. Different mutations in SARS-CoV-2 associate with severe

421         and mild outcome. Int J Antimicrob Agents 2021;57(2):106272.

422    22. Chatterjee P. Covid-19: India authorises Sputnik V vaccine as cases soar to more than 180

423         000 a day. BMJ 2021;373:n978.

424    23. Harvey WT, Carabelli AM, Jackson B, Gupta RK, Thomson EC, Harrison EM, *et al*.

425         SARS-CoV-2 variants, spike mutations and immune escape. Nat Rev Microbiol.

426         2021;19(7):409-424.

427    24. Banoun H. Evolution of SARS-CoV-2: Review of Mutations, Role of the Host Immune

428         System. Nephron 2021;145(4):392-403.

429    25. Watanabe Y, Berndsen ZT, Raghwani J, Seabright GE, Allen JD, Pybus OG, *et al*.

430         Vulnerabilities in coronavirus glycan shields despite extensive glycosylation. Nat

431         Commun 2020;11(1):2688.

432    26. Sanda M, Morrison L, Goldman R. N- and O-Glycosylation of the SARS-CoV-2 Spike

433         protein. Anal Chem 2021;93(4):2003-9.

434    27. vanDorp L, Acman M, Richard D, Shaw LP, Ford CE, Ormond L, *et al*. Emergence of

435        genomic diversity and recurrent mutations in SARS-CoV-2. Infect Genet Evol

436        2020;83:104351.

437    28. Velazquez-Salinas L, Zarate S, Eberl S, Gladue DP, Novella I, Borca MV. Positive

438        selection of ORF1ab, ORF3a, and ORF8 genes drives the early evolutionary trends of

439        SARS-CoV-2 during the 2020 COVID-19 pandemic. Front Microbiol 2020;11:550674.

440    29. Motayo BO, Oluwasemowo OO, Olusola BA, Akinduti PA, Arege OT, Obafemi YD, *et*

441        *al*. Evolution and genetic diversity of SARS-CoV-2 in Africa using whole genome

442        sequences. Int J Infect Dis 2021;103:282-7.

443    30. Pereson MJ, Flichman DM, Martínez AP, Baré P, Garcia GH, Di Lello FA. Evolutionary

444        analysis of SARS-CoV-2 spike protein for its different clades. J Med Virol

445        2021;93(5):3000-6.

**Legends for Figures**

**Figure 1.** SARS-CoV-2 clade distribution pattern in India. **(a)** Pie chart showing the proportion of various clades of the genomes deposited from India in GISAID; **(b)** Schematic geographical map showing the proportion and distribution of clades from different states of India; **(c)** Month wise clade distribution during the year 2020.

**Figure 2.** Phylogenetic tree showing clade diversity for SARS-CoV-2 Indian isolates. These isolates fall under 7 genetic clades with the majority falling under GR clade.

**Figure 3.** Amino acid mutations and their frequency in different regions of S proteins of SARS-CoV-2 isolates from India. Signal peptide (SP), N-terminal domain (NTD), Receptor binding domain (RBD), Protease cleavage site (PC), Fusion peptide (FP), Heptad repeat 1(HR1), Heptad repeat 2 (HR2), Transmembrane domain (TM), Cytoplasm domain (CT).

**Figure 4.** Distribution and frequency of the most prevalent mutations of S protein of SARS-CoV-2 isolates circulated in India during the year 2020. D614G is predominant throughout the year with high frequency followed by L54F mutation. D614G, Q677H and P681H mutations originated during the first half of the year and their appearance was observed throughout the year; L54F, K77M and P812L mutations emerged during the first half of the year but absent after few months of their appearance.

**Figure 5.** Frequency of amino acid mutations impacting O- and N-glycosylation patterns. Few sites such as S221, T602 and N1074 are having more than one mutation.
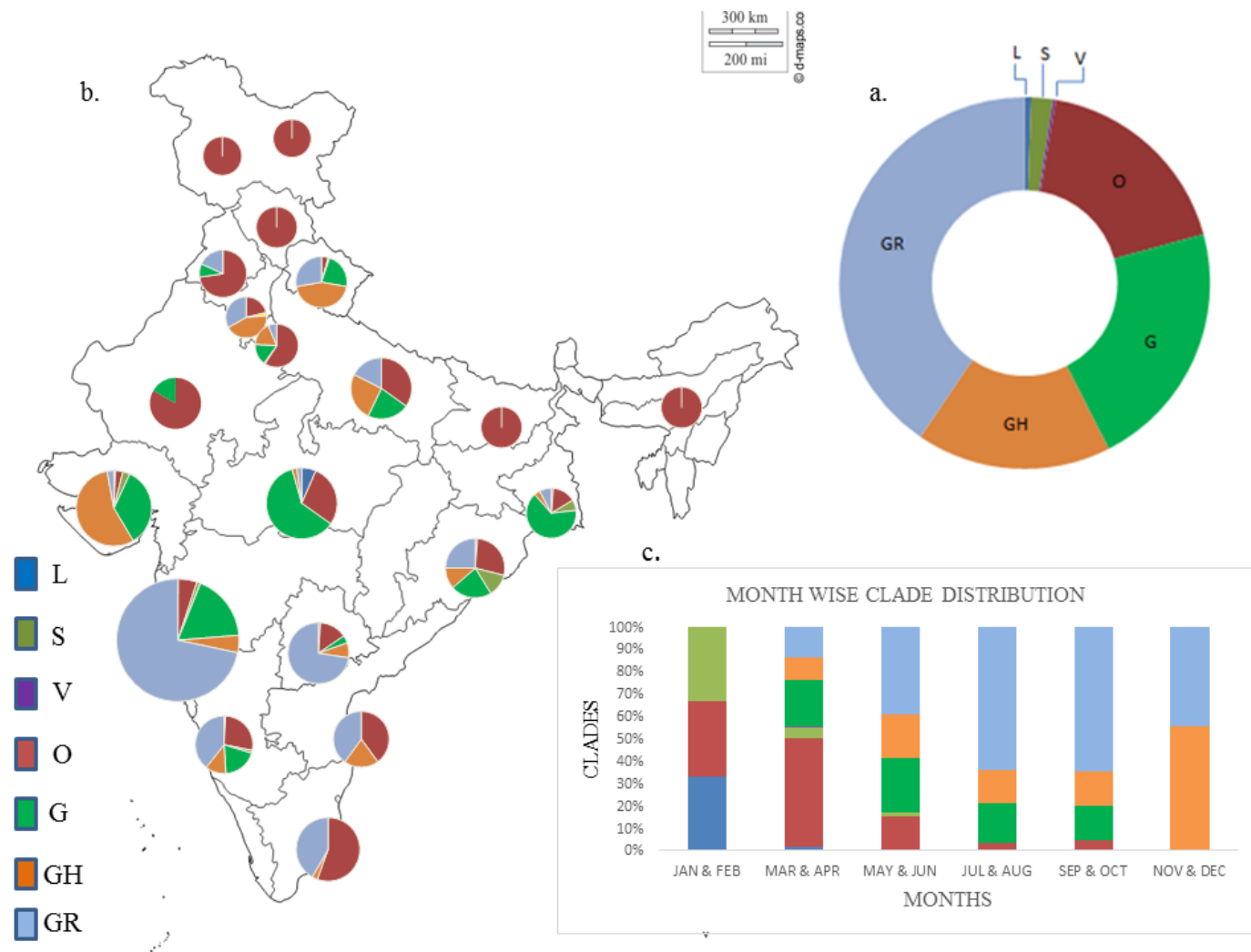
**Figure 6.** Phylogenetic tree of isolates having distinct mutations in the gene of S protein. The tree was constructed by Maximum-Likelihood method with the tree having the root as SARS-CoV-2 Wuhan-Hu-1 sequence (NC_045512.2).
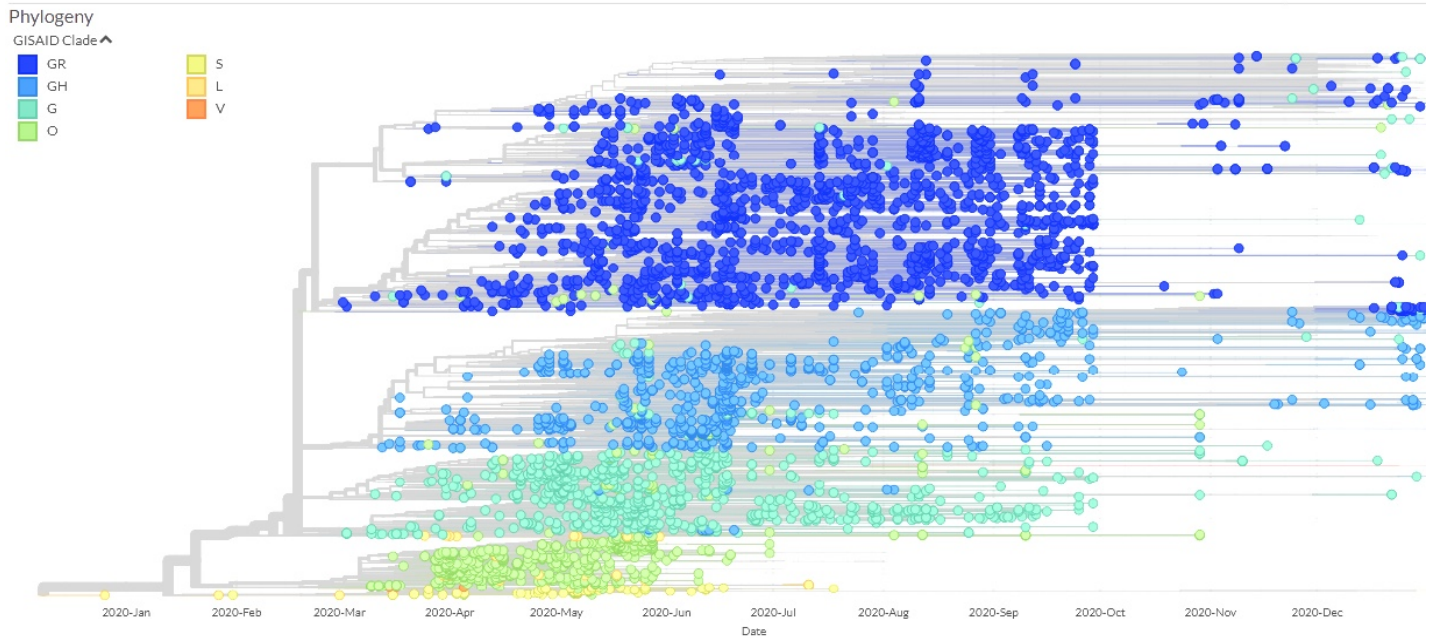
**Legends for Tables**

**Table 1.** Aminoacid substitution mutations observed across various regions of S proteins of Indian SARS-CoV-2 isolates.

**Table 2.** N- and O-linked glycosylation sites of S protein of SARS-CoV-2 and amino acid mutations at these sites affecting the glycosylation pattern in Indian SARS-CoV-2 variants.
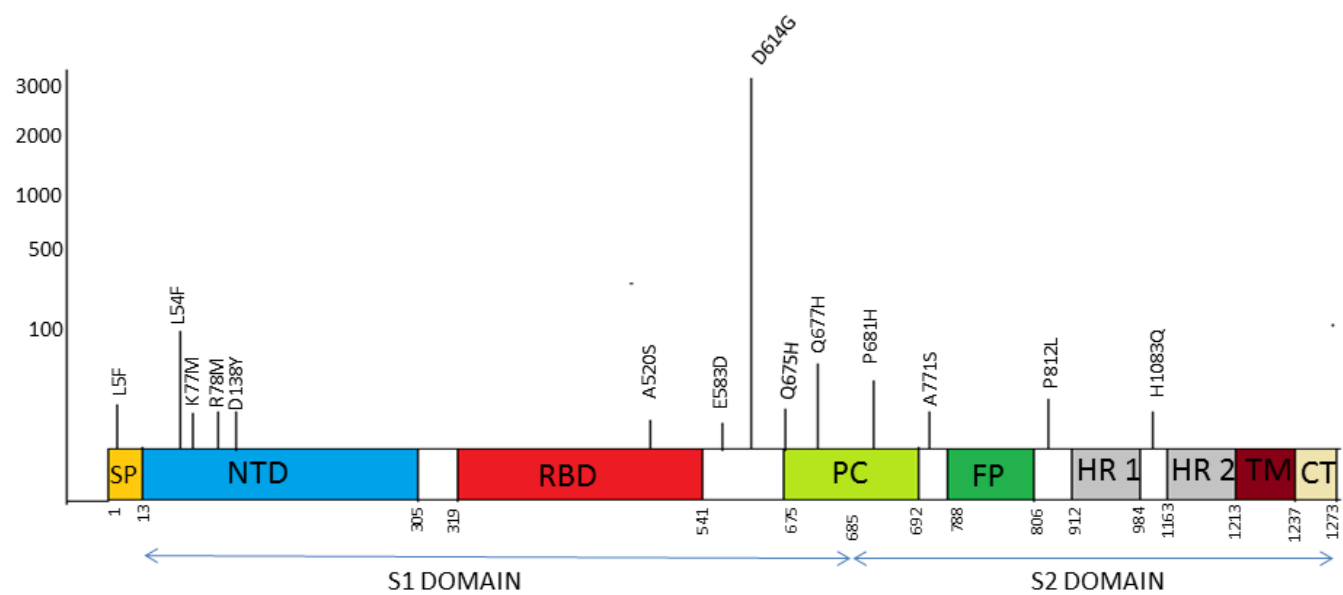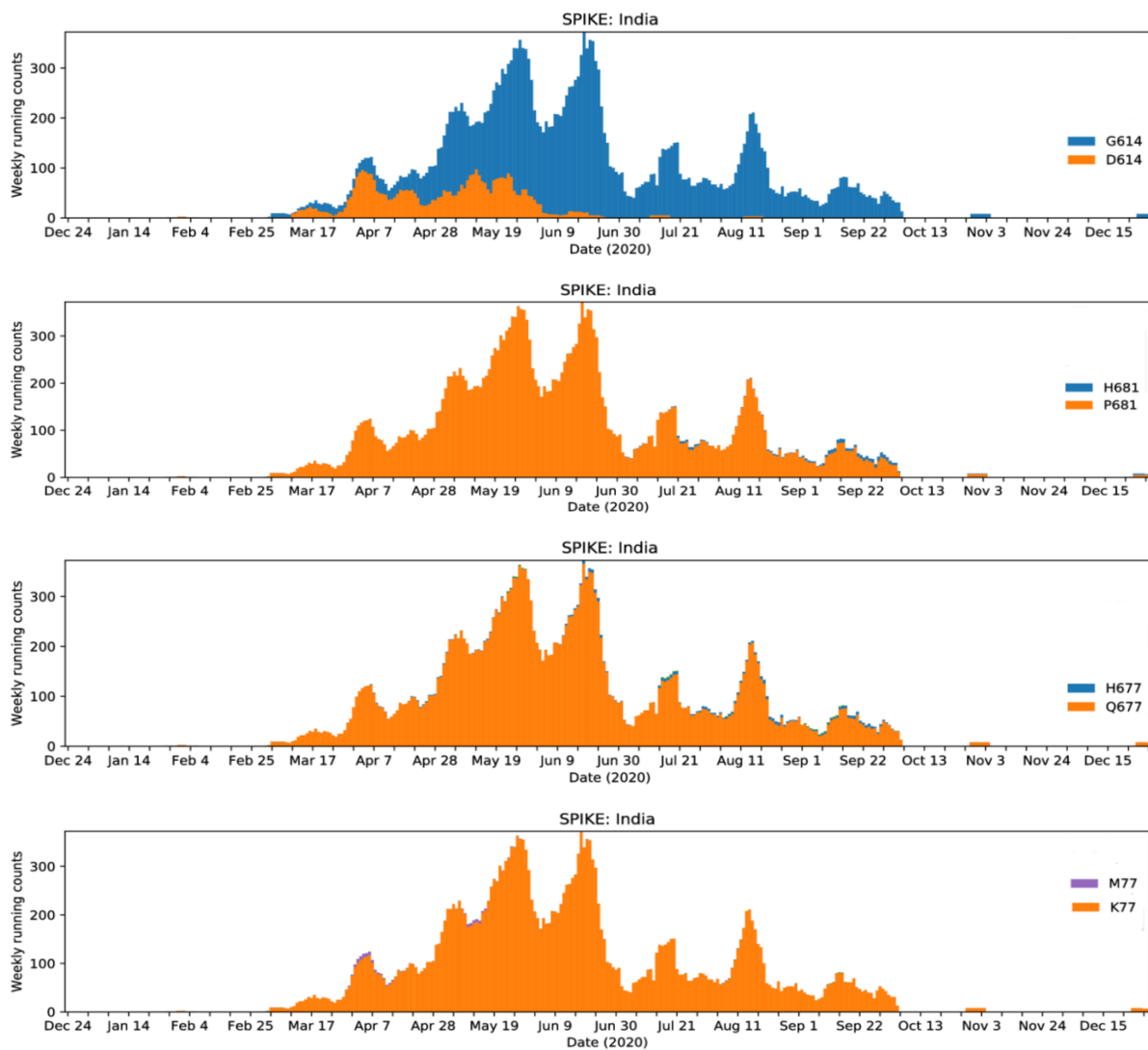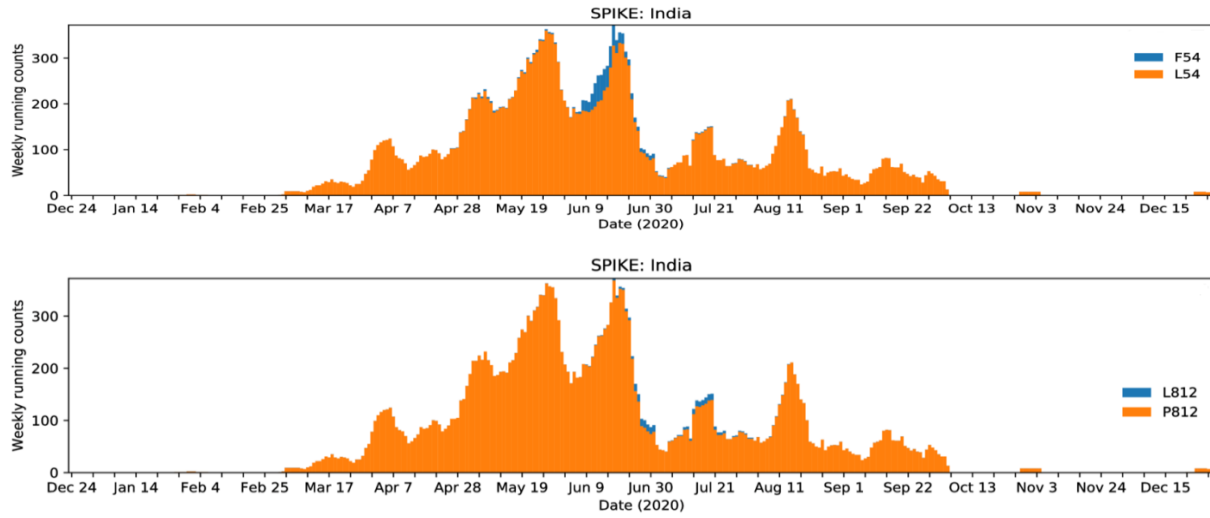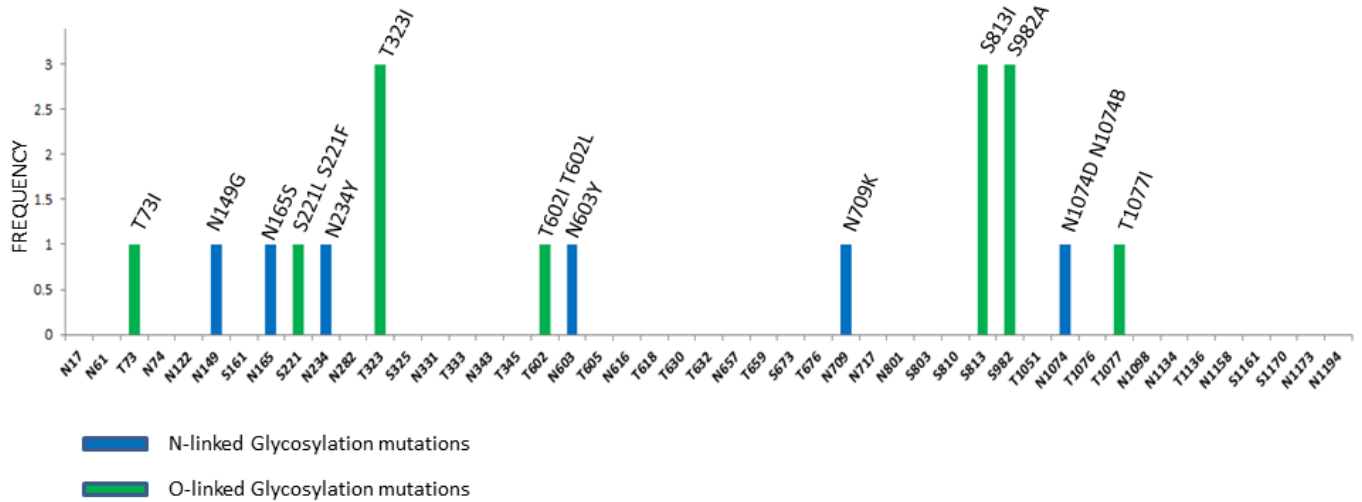
**Figure 1**

# Figure 2

**Figure 3**

**Figure 4**

**Figure 5**



N - linked and O - linked Glycosylation sites in spike protein

**Figure 6**

| Region | Position | Number of mutation sites | Number of mutations |
|--------|----------|--------------------------|---------------------|

hCₒ

hCoV-1ₐ

hCoV

hCoV

hCoV-19/India

hCoV-19/Ind

hCo

ᵣ

**Table 1**

| | | | |
|---|---|---|---|
| Signal Peptide | 1 - 13 | 7 | 9 |
| N-Terminal Domain | 14 - 305 | 144 | 211 |
| Receptor Binding Domain | 319 - 541 | 53 | 63 |
| Protease Cleavage Site | 675 - 692 | 8 | 11 |
| Fusion Peptide | 788 - 806 | 6 | 6 |
| HR1 | 912 - 984 | 24 | 31 |
| HR2 | 1163 - 1213 | 15 | 18 |
| Transmembrane Domain | 1214 - 1237 | 13 | 16 |
| Cytoplasm Domain | 1238 - 1273 | 12 | 13 |

**Table 2**

| N-linked | Mutation | O-linked | Mutation |
|---|---|---|---|

| Glycosylation Site (NGS) | | Glycosylation Site (OGS) | |
|---|---|---|---|
| N17 | - | T73 | T73I |
| N61 | - | S161 | - |
| N74 | - | S221 | S221, S221F |
| N122 | - | T323 | T323I |
| N149 | N149G | S325 | - |
| N165 | N165S | T333 | - |
| N234 | N234Y | T345 | - |
| N282 | - | T602 | T602I, T602L |
| N331 | - | T605 | - |
| N343 | - | T618 | - |
| N603 | N603Y | T630 | - |
| N616 | - | T632 | - |
| N657 | - | T659 | - |
| N709 | N709K | S673 | - |
| N717 | - | T676 | - |
| N801 | - | S803 | - |
| N1074 | N1074D, N1074B | S810 | - |
| N1098 | - | S813 | S813I |
| N1134 | - | S982 | S982A |
| N1158 | - | T1051 | - |
| N1173 | - | T1076 | - |
| N1194 | - | T1077 | T1077I |
| | | T1136 | - |
| | | S1161 | - |
| | | S1170 | - |
| | | S1175 | - |