

1 **Unique mutational changes in SARS-CoV2 genome of different state of** 2 **India**

3 Priti Prasad^{1,2*}, Shantanu Prakash^{3*}, Kishan Sahu^{1,2}, Babita Singh¹, Suruchi Shukla³, Hricha Mishra³, Danish
4 Nasar Khan³, Om Prakash³, MLB Bhatt³, SK Barik¹, Mehar H.Asif^{1,2}, Samir V. Sawant^{1,2}, Amita Jain^{3#}, Sumit
5 Kr. Bag^{1,2#}

6

7 ***Equal Contribution**

8

9 ¹CSIR-National Botanical Research Institute, 436, Rana Pratap Marg, Lucknow - 226001, India

10 ²Acadaemy of Scientific and Innovative Industrial Research (AcSIR), Ghaziabad -201002, India

11 ³King Georges Medical University, Lucknow-226001, India

12

13 **#Corresponding author**

14

15 **Amita Jain**

16 **Professor and Head, Department of Microbiology**

17 **King Georges Medical University**

18 **Lucknow-226001, India**

19 **Email: amita602002@yahoo.com**

20

21 **Sumit Kr. Bag**

22 **Principal Scientist, Computational Biology**

23 **CSIR-National Botanical Research Institute**

24 **Lucknow-226001, India**

25 **Email: sumit.bag@nbri.res.in**

26

27

28

29

30

31

32

33

34

35

36

37 **Abstract**

38 COVID-19 leads to a global emergency that causes more than 7 million casualties until mid-
39 August throughout the world. In India alone, 2 million confirmed cases were reported that
40 increased abruptly day by day with the lowest fatality rate. The availability of a large number
41 of Severe Acute Respiratory Syndrome Coronavirus-2 (SARS-CoV2) genome in the public
42 domain provides a great opportunity to study mutational changes in COVID-19 genomes in
43 Indian populations. In this study, we sequenced the genomes of SARS-CoV2 viruses isolated
44 from 47 individuals from 13 districts of Uttar Pradesh (UP), the largest state of India using
45 Third Generation Sequencing Technology. We further did the phylogenetic clustering of UP
46 state of Indian SARS-CoV2 genomes revealed a perceivable point that no UP samples were
47 aligned on the USA defined clade where the fatality rate is high. We also identified 56
48 distinctive Single Nucleotide Polymorphism variations in UP state that majorly clustered into
49 two groups which shows the deleterious effects on the genome. Additionally, we conducted
50 the mutation analysis of the 2323 SARS-CoV2 genome of different states of India from the
51 Global Initiative on Sharing All Influenza Data (GISAID) where we find ~80% unique
52 mutations rate in each sample of the Indian population. Thus, this is the first extensive
53 mutational study of the largest state of Indian populations in which we report the novel
54 deleterious SNPs in virus genome along with the other states which access the less infectious
55 form of SARS-CoV2 genome through synonymous to nonsynonymous mutation variation.

56 **Keywords**

57 Third Generation Sequencing, COVID-19, India, Uttar Pradesh, SNP, Unique

58

59

60 **Introduction**

61 The Coronavirus disease (COVID-19) emerged out as a global pandemic[1] and has become
62 a center point for many researchers to identify the potential region/s of its pathogenicity. The
63 causative agent, SARS-CoV2 belongs to the Coronaviridae family comprises of single-
64 stranded, positive-sense RNA with a genome size of approximately 30 kb [2]. SARS-CoV2
65 genome consists of 29 open reading frames (ORFs), among which four are structural proteins
66 (envelope protein, membrane protein, nucleocapsid protein, and Spike protein), some
67 accessory protein translated by ORF3a, ORF6, ORF7a, ORF7b, ORF8, and ORF10 and 16
68 nonstructural proteins (nsps) enclosed the complete genome [3].

69 At the end of December 2019, China reported its first case in Wuhan city, after then it grew
70 exponentially through human to human transmission. During its spreading, it evolved
71 continuously which causes changes in fatality rates in different countries [4][5]. As of 20th
72 August, approximately 23 million confirmed cases were reported worldwide with ~800
73 thousand registered deaths. At present India is ranked the 3rd position in the total number of
74 cases, however, the fatality rate is declining to 1.9 %, which is considered as the one of the
75 lowest fatality rates [6] [7] [8]. Many reports correlate the fatality rate to the mutational
76 changes in the virus genome in different geographical regions [9][10], [11][12] in which
77 D614G substitution at the spike protein cause a higher mortality rate in human [4]. A study
78 on western Indian populations linked the mutational changes to the specific age group and
79 their mode of infection. To delineate the different mutations rate and understanding the
80 evolution in the SARS-CoV2 genome of different continents, many genomes were sequenced
81 by employing different techniques and are submitted at GISAID [13]. More than 84
82 thousands genome sequences of the full-length SARS-CoV2 were available on the GISAID
83 platform where India has sequenced >2400 viral genomes and from where Indian specific
84 clade was retrieved [14].

85 India comprises of 29 states with the 7 Union Territories in which Uttar Pradesh is ranked 5th
86 in confirmed cases and surprisingly, only 11 sequenced samples were available on the
87 database. Therefore in this study, we first report the detailed genome sequences of SARS-
88 CoV2 in different districts of Uttar Pradesh (UP) through the third generation sequencing
89 technology (Oxford Nanopore PromethION) and generate the phylogenetic clustering . We
90 have also investigated the unique mutational features in a different state of India in
91 comparison to other continents of the world to find out the mutation rate during the
92 transmission. We also correlated the fatality rate of different state with their synonymous to
93 nonsynonymous mutation ratio.

94 **Methods**

95 **RNA Preparation and Sequencing of the viral genome**

96 Total viral nucleic acid was extracted from the clinical sample of SARS-CoV-2 using
97 PureLink DNA/RNA mini kit (Invitrogen, USA) as per manufacturer's instructions and
98 Indian Council of Medical Research (ICMR) guidelines [15] . The Human RnaseP gene was
99 tested in all samples to check the quality of samples and to validate the process of nucleic
100 acid extraction. Real time RT-PCR assays were performed using Superscript III One-Step
101 Real-Time PCR kit (ThermoFisher Scientific, USA) as per the protocol[16]. All samples
102 were tested by 2 stage protocol consisting of screening and confirmation. E-gene and RnaseP
103 were used for screening while RdRp and ORF1b were used as a confirmatory assay for
104 SARS-CoV-2. All the samples positive for E-gene, RdRp, and ORF1ab with Ct \leq 25 were
105 further taken for whole-genome sequencing.

106 Virus genomes were generated by using ARTIC COVID-19 multiplex PCR primers followed
107 by nanopore sequencing on two ONT PromethION flow cells[17]. For PromethION
108 sequencing, libraries were prepared using the ligation sequencing kit (SQK-LSK109) and

109 Native barcoding (EXP-NBD 104 and EXP-NBD 114) of 24 samples in a single flow cell.
110 The selected samples were converted to first-strand cDNA using reverse transcriptase as
111 suggested in Nanopore protocol “PCR tiling of COVID-19 virus” (Version:
112 PTC_9096_v109_revD_06Feb2020). The multiplex PCR was performed with two pooled
113 primer mixture as recommended and the cDNA was amplified. After 35 rounds of
114 amplification, the PCR products were collected, purified with AMPure XP beads (Beckman
115 coulter), and quantified using Qubit4 Fluorometer (Invitrogen). The double-stranded cDNA
116 was further subjected to end-repair, A-tailing, and Native barcode ligation individually. The
117 barcoded samples were pooled together and ligated with sequencing adapters followed by
118 loading of 50ng of pooled barcoded material and sequencing on 2 PromethION flow cells
119 having 24 samples each.

120 **Genome Assembly -**

121 Quality passed nanopore reads were subjected to remove the chimeric reads between the 300
122 to 500 base pair length according to the amplicon size. With the help of the Artic network
123 pipeline [18], we individually assembled the high coverage (approx. 9000X) filtered long
124 reads of each sample. During the assembly --normalize options were used to normalizes the
125 data up to 200 bp coverage. This normalization is useful for the minimization of the
126 computational requirements of our high coverage data. The SARS-CoV2 reference genome
127 (MN908947.3) [19] was used for generating the consensus sequence from the assembled
128 reads. Low complexity or primer binding regions were masked during the consensus
129 sequence generation and variants calling to avoid the ambiguity.

130 **Data Preparation-**

131 For comparative and explanatory analyses of our assembled genome form Uttar Pradesh (UP)
132 district of India, we downloaded the whole genome data of all Indian samples from the

133 GISAID database that was available until the 12th of August. As an outlier, and take the
134 representation of all the continents of the world, we randomly selected 100 samples from
135 Asia, Europe, South America, North America, Central America, and Eurasia. A total of 2923
136 processed samples (2323 from India and 600 from different continents) were further screened
137 based on unrecognized nucleotides (Ns) with the full length of the genome. We kept only
138 those samples whose unrecognized base pair number (number of Ns) is less than 1 % of the
139 total sequence followed by the masking of the low complexity region.

140 **Phylogenetics clustering of Indian sample –**

141 We sequenced and assembled a total of 47 genome sequences from the Uttar Pradesh state of
142 India followed by the phylodynamic clustering according to the nextstrain clade [20]. Whole-
143 genome alignment has been carried out by the mafft aligner [21] by taking the Wuhan
144 nCov2019 Hu-1(MN908947.3) sample as a reference genome. IQTREE software [22] was
145 used for the phylogenetic tree with the GTR substitution model. The refinement of the tree
146 was done by the --timetree parameters to adjust the branch lengths and filter out the
147 interquartile ranges from regression. Thus, the refined time-resolved tree was further
148 screened to derive the ancestral traits, infer nucleotide mutation as well as an amino acid
149 variation from the root. The resulting phylodynamics tree was then used to defining the clade
150 based on a new classification that signifies the different geographic spreading of nCoV-2019
151 in India. The final tree was further visualized through the nextstrain shell by using the auspice
152 interface [23].

153 **Divergence estimates of the newly assembled genome-**

154 Divergence estimation of the newly assembled genome and the most recent common ancestor
155 (tMRCA) was examined using the BEAST v1.10.4 program [24]. Through the MRCA
156 analysis, the origin of the ancestral virus has been identified that was present in our samples.

157 With the strict molecular clock and the exponential growth model, we processed the Markov
158 chain Monte Carlo algorithm on the masked alignment file for 100 million steps, where the
159 first 30 million steps were used for burn-in to get the effective sample size. Bayesian
160 coalescent analysis of the UP samples was conducted using the GTR substitution model with
161 the kappa scale of 1 to 1.25. Molecular clock and root divergence time were reconfirmed
162 with the treetime [25] using the same substitution model.

163 **Variant identification and functional evaluation**

164 We scrutinized the unique mutations in our assembled genome of UP district that was called
165 during the assembly. The SNP variants between all processed samples (downloaded and
166 assembled) and to ensure the uniqueness of variants in our assembled data, we aligned each
167 downloaded samples separately to the reference genome (MN908947.3) by use of mafft
168 aligner [21] and find the variation with the snp-sites program and [26] nextstrain. The
169 resulting Variant Call Format (VCF) file was annotated further to find out the position of the
170 mutation that affects the nCoV-2019 genome though the snpEff program [27]. Finally,
171 manual screening of mutation has been carried out to define the unique mutation in our
172 processed samples (2370) along with the other continents of the world. Functional validation
173 of the emerging new mutations in UP samples was also assessed through the Sorting
174 Intolerant from Tolerant (SIFT) [28]. SIFT evaluated the functional consequences of the
175 variants where the SIFT score of 0.0 to 0.5 is used to explain the deleterious effect while
176 >0.05 was interpreted as a toleratable mutation. Visualization of identified variants of Uttar
177 Pradesh samples was processed through the R package [29].

178 **Results and Discussion:-**

179 **Demographic Representation**

180 In this study, a total of 47 genomes of SARS-CoV2 have been sequenced using the long reads
181 PromethION sequencer. All the samples are from the Uttar Pradesh (UP) state, India
182 encompassing over 13 different districts (**Fig 1a**). The samples were collected from
183 symptomatic and asymptomatic patients in which the male-female ratio of patients was 3:1.
184 The age of the patients varied from 2 to 86 years with an average age of 35years. The
185 maximum number of samples was collected in May month while 2 and 13 samples were
186 collected in March and April respectively (**Supplementary File 1**). The average coverage of
187 all the samples was 9000X which is correlated with the ct value of the samples (**Fig. 1b**).

188 **Assembly of the whole genome**

189 The resulting fastq and fast5 files of the nanopore PromethION sequencer were quality
190 checked by the artic network pipeline [18]. The amplicon reads of 300 to 700 base pair length
191 was removed from the raw fastq file to avoid the generation of chimeric reads. Filtered fastq
192 file and raw fast5 format files of each sample were taken as input for whole-genome
193 assembly. The MN908947.3 genome was used as a reference to generated the consensus
194 sequence for each sample individually. Additionally, medaka parameters were also used to
195 identify the variants or SNP based on assembled reads during the generation of the consensus
196 sequence. To avoid the ambiguity at the time of SNP calling, the primer binding region in the
197 genome was masked. The resulting Single-nucleotide Polymorphism (SNPs) were considered
198 as true SNPs in our sequenced samples data because they have high-quality read depth (~300)
199 at each base pair. These High-quality SNP reads were used to define the unique mutation
200 features in the SARS-CoV2 genome samples of UP state in comparison to the samples from
201 other states of India as well as from different continents. The assembled whole-genome
202 sequence was submitted to the Global Initiative on Sharing All Influenza Data with accession
203 ID (EPI_ISL_516940-EPI_ISL_516986) and their metadata information is available in
204 **Supplementary File 1**.

205 **Phylogenetic analyses of the assembled genome**

206 Phylogenetic clustering is used for clade assignment in our assembled genome based on
207 nextstrain where we can trace the origin of the clade concerning the specific mutation on the
208 genomic region. A total of 47 assembled genome of SARS CoV2 from the 13 different
209 districts of UP in India and 856 genomes from 21 different states of India (available on 10th
210 June) was used to determine the phylogenetics clustering according to the new clade
211 classification by taking the Wuhan SARS-CoV2 genome reference genome
212 (EPI_ISL_406798) as an outlier. All the information regarding the samples including
213 accession number, submitting institution name, collection date is mentioned in

214 **Supplementary File 2.**

215 The augur pipeline [30] was used for the alignment of all the samples, their phylogenetic
216 clustering, and deciphering the time resolving tree, using SARS-CoV2 (MN908947.3)
217 genome as the root. The resulting tree was subjected to defining the potential emerging clades
218 based on nucleotide mutation(s). Indian samples were clustered into all 5 clusters namely
219 19A, 19B, 20A, 20 B, 20C (**Fig. 2a**) as it is represented in the nextstrain nCoV global page
220 and described in [31]. The two clades 19A and 19B are primitive ones that are dominated by
221 the Asia continents. These two clades are delineated by C to T mutation at 8782 positions at
222 ORF1ab region and T to C mutation at 28144 positions. A total of 43% of our processed
223 sample aligned on these two clades of which samples from the UP states aligned only on
224 clade 19A.

225 The clade 20A, 20B, and 20C are well known for the European and North American clades
226 respectively, which deviated from the root (19A) in early 2020. Most of the Indian samples
227 (around 41%) were established under 20A clade with C14408T and A23403G mutations in
228 the orf1ab region followed by the 19A clade were 332 samples were superimposed (**Fig. 2b**).

229 The smallest cluster of India samples has been seen in the case of a 20C clade. The USA is
230 the primary source for 20C clade and only 4 Indian samples (2 from Delhi and Gujarat each)
231 were aligned. This 20C clade formed by Delhi samples (EPI_ISL_435063, EPI_ISL_435064)
232 doesn't show any nucleotide or amino-acid mutation whereas Gujarat samples
233 (EPI_ISL_435051, EPI_ISL_435052) recorded ORF1ab variations (mention the variations).
234 The SARS-CoV2 genome samples of UP state doesn't align on the 19B and 20C clade which
235 is a matter of investigation that can be solved with the increasing number of COVID genome
236 sequences of Uttar Pradesh state with the probable mutation of clade defining mutation like
237 1059T for 19B and 25563T for the 20C clade.

238 **Estimation of the date of origin**

239 The mutation rate of 47 sequenced genome in UP state was estimated using BEAST software
240 and by using the coalescent model as the prior of the trees [24]. The estimated rate of
241 mutation is predicted as $3.279e-04$ substitutions per site per year for the given SARS-CoV2
242 genomes samples. The 95% of Highest Posterior Density intervals (HPD) are in the range of
243 $2.24e-04$ to $4.459e-04$, which was calculated for each sample (**Supplementary File 3**). The
244 substitution rate is further confirmed through the treetime software [25] where root to tip
245 regression rate is calculated as $3.078e-04$ with the 0.05 correlation value (**Fig. 3**). This
246 correlation value indicates the informative behavior of inputs for temporal information to
247 rationalize the molecular clock approach. The root date of the SARS-CoV2 genome from UP
248 states is evaluated as 29 December 2019 which substantiated the reported timing of the
249 SARS-CoV2 genome in Wuhan city of China [32].

250 **Assessment of unique Mutational features in different states of India**

251 We identified the mutational features in the SARS-CoV2 genome of other states of India also
252 wherein the unique mutation rate is around 80% of the total sequenced genome. The unique

253 mutation features hypothesize that at each transmission, the virus mutated at the higher speed
254 with insynonymous to nonsynonymous mutation rate is around 0.5. For the isolation of
255 unique SNPs features, we selected only those states whose SARS-CoV2 genome samples
256 were available in more than 10 samples. Thus, total 12 states; Gujarat, Maharashtra,
257 Telangana, Delhi, Odisha, Karnataka, West Bengal, Uttarakhand, Madhya Pradesh, Tamil
258 Nadu, Haryana, and Uttar Pradesh were selected (**Table 1**) whose average active case is
259 around 2 lacs according to the WHO India tally [7]. We strike off the common mutation
260 position from different continents (Asia, Africa, Eurasia, North America, Oceania, and South
261 America) to check the mutated version of the SARS-CoV2 genome in the Indian strain. A
262 higher number of unique mutations were seen in Maharashtra followed by the Gujarat,
263 Odisha, Telangana, West Bengal, and Karnataka where an abundant amount of sequenced
264 data were also available (**Table 1**). These unique SNPs information were annotated on the
265 SARS-CoV2 genome and classified it into synonymous, missense, stop gain, stop loss,
266 upstream gene variant and downstream gene variant (**Fig. 4**). Surprisingly, Delhi has very
267 low number of unique mutation rate (23%) and its mortality rate is comparably high (2.66)
268 which suggest that Delhi was infected with a fatal version of the virus genome. The
269 phylogenetic clustering by nextstrain also supplemented the information where genome from
270 Delhi was clustered on 20C clade [20] which is specifically defined as the North American
271 clade. On the other hand, Maharashtra and Gujarat were top listed according to the fatality
272 rate where the unique mutation rate is ~111% and ~89% respectively. These two states have
273 the equal number of missense variant which affect almost all the genomic region with the
274 stop codon gain at the envelope protein in Maharashtra state and ORF7a in Gujarat state
275 (**Supplementary File 4**) (**Fig. 4**). Odisha state has also higher number of mutation rate with
276 the synonymous to nonsynonymous mutation rate is below (0.39) than average (0.5) and its
277 fatality rate is very low. Similar pattern is observed in the Telangana state as well as in

278 Harayana state so through this we can postulated that fatality rate could be estimatd through
279 the synonymous to the nonsynonymous ratio of unique mutational features.

280 **Dynamic characterization of Unique SNP variation in UP state**

281 The identification of clade is defined by at least two mutations changes from its parent clade
282 which is ubiquitous for all SARS-CoV2 genome. Besides the mutations that are shaping the
283 clade, we identified 56 distinctive SNPs in our UP district sequence samples, referred to as
284 unique SNPs. Among these uniquely identified SNPs, 36 SNPs are annotated as a missense
285 variant which causes to change in the amino acids while 2 mutations (T18126A, C18431T)
286 show the premature stop gain in the ORF1ab regions (**Table 2, Fig 5**). The functional
287 evaluation of the nonsynonymous mutations was investigated through Sorting Intolerant
288 From Tolerant SIFT) prediction server which revealed the deleterious effects of 37% of total
289 mutations. A Maximum number of deleterious mutations have been seen in the non-
290 structural protein of ORF1ab regions while 2 stop gain mutations have also been seen in the
291 NSP11 margin [2]. The Threonine to Isoleucine intolerable mutations in ORF1ab regions, at
292 6056th position, affects the NSP11 protein in the NBRI-N21 genome which might change the
293 polymerase activity of the genomes. The different type of identified mutation in non-
294 structural protein 11(NSP11) of orf1ab in UP states gives a clue for conducting the detailed
295 study on it, as it is associated with the programmed cell death evasion, cell cycle and DNA
296 replication [35][36]. Only 1 each deleterious mutation (**Table 3**) was found in ORF3a,
297 ORF7a, and ORF8 regions respectively which functions as an accessory factor in the viral
298 genome [37] and also assisted to virus-host interactions [38]. Accessory factor, ORF3a
299 regulates the interferon signaling pathway and cytokines production [39] and thus is involved
300 in the virus pathogenesis. Such kind of deleterious mutation in the accessory protein has been
301 found in the latest evolving SARS-CoV2 genome from Indian samples [40]. A structural
302 protein, N encoded the nucleocapsid protein shows the deleterious effects by the mutations at

303 D371V, K370N, D371Y, P279Q, and T362I positions in 4 UP state genomes. SARS CoV2 N
304 proteins have a highly conserved domain [41] that plays a crucial function to complete the
305 viral life cycle by regulating the RNA transcription, replication, and modulating the infected
306 cell biological processes [42]. Another nonsynonymous, deleterious mutation is noticeable
307 in NSP3 protein in the NBRI-N42 sample. This protein has papain-like viral protease activity
308 to generate the other replicase subunits from nsp1 to nsp16 [43]. This intolerable mutation
309 from Valine to Leucine at 1570 position might affect the functionality of the papain activity
310 which leads to the inactivation of replicase subunits from nsp1 to nsp16.

311 **Conclusion:**

312 The present study is the first one to report the SARS-CoV2 genome sequence of the Uttar
313 Pradesh state with a large scale of high coverage third-generation sequencing data. We
314 identified 56 peculiar SNPs in our sequenced genome in which more than 40% of SNPs in the
315 ORF1ab region showed the deleterious effects. These SNP variations in the ORF1ab region
316 might affect the replication of the virus genome during the infection which ultimately would
317 be the reason for the less fatality rate. The identified mutations in SARS-CoV2 genomes of
318 47 individuals could be used as the potential target for personalized medications or effective
319 vaccine doses to combat the effects of the COVID-19. Additionally, the relation with the
320 synonymous to nonsynonymous unique mutation ratio with the fatality rate could be studied
321 further to understand the putative region/SNPs that cause fatal to the human being.

322 **Author contributions:**

323 Priti Prasad: Conceptualization; Data curation; Supervision; Formal analysis; Investigation;
324 Methodology; Resources; Software; Visualization; Writing - original draft; Writing - review
325 & editing, Shantanu Prakash: Conceptualization; Methodology; Writing - review & editing,
326 Mehar H. Asif: Conceptualization; Supervision; Writing - original draft; Writing - review &

327 editing, Kishan Sahu: Resources; Writing - review & editing, Suruchi Shukla, Babita Singh,
328 Hricha Mishra, Danish Nasar Khan, Om Prakash, MLB Bhatt: Resources, SK Barik:
329 Conceptualization; Resources; Writing - review & editing, Samir V. Sawant:
330 Conceptualization; Resources; Supervision; Formal analysis; Writing - review & editing,
331 Amita Jain: Conceptualization; Resources; Writing - review & editing, Sumit Kr. Bag:
332 Conceptualization; Resources; Supervision; Writing - review & editing

333 **Conflicts of Interest**

334 Authors declare no conflicts of interest

335 **Acknowledgments**

336 Authors acknowledge the GISAID team and all those who submitted the genome to the
337 GISAID database without which it's impossible to conduct the research. PP and KS
338 acknowledged the University of Grant Commission for providing the Senior Research
339 fellowship. Institute manuscript number is CSIR-NBRI_MS/2020/08/04.

340 **Figure legends:**

341 **Fig. 1**

342 a) The detailed map of the SARS-CoV2 genome sequence of different districts of the Uttar-
343 Pradesh state of India. The scale denoted the number of genome sequencing of the virus
344 genome. Grey color means no genome was sequenced while the darker blue color represents
345 the highest number of genome sequencing. b) The Co-linear plot of the Ct value of 47
346 positive COVID-19 patients to the coverage of the whole genome sequence long-read data.
347 Each sample was represented through the blue dots where higher Ct value shows high
348 coverage data.

349 **Fig. 2**

350 a) The phylogenetic clustering of 889 samples of India into 5 different clades according to the
351 nextstrain based mutation. The 19A and 19B clades are the primitive clades that were
352 sequenced in early 2020; whereas 20A, 20B and 20C clades are emerging clades. The
353 reference genome of the Wuhan sample was implemented in the clustering trees which
354 aligned in 19A clade before the starting of the Jan-2020. b) The distribution of genome
355 sequenced in different clades was visualized through the piechart. In the 20A clade, 366
356 genomes were clustered followed by the 19A clade. The lowest number of clustering has
357 been seen in the 20C clade which is well known for the USA defining clade.

358 **Fig. 3**

359 The root-to-tip regression rate of sequenced SARS-CoV2 genome of Uttar Pradesh state of
360 India to investigate the origin of infection. X-axis describes the timing of the sample
361 collection of COVID-19 infected persons. Y-axis describes the regression rate in comparison
362 to the sample collection date. The origin of infection in Uttar Pradesh samples is estimated by
363 2019.5.

364 **Fig. 4**

365 Unique mutation features in 12 different states of India whose SARS-CoV2 genome sequence
366 is present in large amounts in GISAID. The pie chart shows the different types of variants
367 that were annotated on the SARS-CoV2 genome with their unique mutation variation per
368 high-quality genome sequenced were available on the public domain. In the pie chart of Uttar
369 Pradesh, the publicly available genome and the assembled genome in this study were
370 mentioned separately.

371

372 **Fig. 5**

373 The nonsynonymous mutation of the SARS-CoV2 genome of the Uttar Pradesh 47 samples.
374 The scale bar shows the length of the genome on which the dashed line with a brown color
375 number represents the mutational points that putatively altered the functionality of the
376 genome. The circled nucleotides are the mutated nucleotide bases, in which black color bases
377 are the missense variation while red color shows the stop gain at the genomic regions.

378 **Table 1**

379 The statistics of the different selected states of India along with their unique mutation rate
380 and the fatality rate for the respective states according to the mid-August, 2020 COVID
381 cases. The sequences data is downloaded from GISAID, available till 12th of August.

382 **Table 2**

383 Detail of Nonsynonymous SNP mutation of all 47 sequenced UP samples with their genomic
384 and amino acid position on the genome. The different colored Gene color shows different
385 genomic regions while stop codon affecting mutation were highlighted in green in color.

386 **Table 3**

387 The Putative nonsynonymous mutation with their SIFT score. The SIFT score value less than
388 0.05 is depicted as a deleterious mutation that is highlighted in yellow whereas unhighlighted
389 SIFT score value for the particular SNPs is predicted as a tolerable mutation.

390 **References -**

- 391 [1] P. Yang and X. Wang, "COVID-19: a new challenge for human beings," *Cell. Mol. Immunol.*, vol. 17,
392 no. 5, pp. 555–557, 2020, doi: 10.1038/s41423-020-0407-x.
- 393 [2] J. F. W. Chan *et al.*, "Genomic characterization of the 2019 novel human-pathogenic coronavirus
394 isolated from a patient with atypical pneumonia after visiting Wuhan," *Emerg. Microbes Infect.*, vol. 9,
395 no. 1, pp. 221–236, 2020, doi: 10.1080/22221751.2020.1719902.

- 396 [3] R. A. Khailany, M. Safdar, and M. Ozaslan, "Genomic characterization of a novel SARS-CoV-2," *Gene*
397 *Reports*, vol. 19, no. April, p. 100682, 2020, doi: 10.1016/j.genrep.2020.100682.
- 398 [4] M. Becerra-Flores and T. Cardozo, "SARS-CoV-2 viral spike G614 mutation exhibits higher case
399 fatality rate," *Int. J. Clin. Pract.*, no. May, pp. 4–7, 2020, doi: 10.1111/ijcp.13525.
- 400 [5] B. Korber *et al.*, "Spike mutation pipeline reveals the emergence of a more transmissible form of SARS-
401 CoV-2," *bioRxiv*, p. 2020.04.29.069054, 2020, doi: 10.1101/2020.04.29.069054.
- 402 [6] H. Ritchie *et al.*, "India : Coronavirus Pandemic Country Profile," pp. 1–24, 2020.
- 403 [7] World Health Organization, "Coronavirus disease," *Coronavirus Dis. Situat. Rep. – 119*, vol. 2019, no.
404 May, p. 2633, 2020, doi: 10.1001/jama.2020.2633.
- 405 [8] F. Welfare, "Latest Updates 27.04.2020," pp. 3–5, 2020.
- 406 [9] D. Mercatelli and F. M. Giorgi, "Geographic and Genomic Distribution of SARS-CoV-2 Mutations,"
407 *Front. Microbiol.*, vol. 11, no. July, pp. 1–13, 2020, doi: 10.3389/fmicb.2020.01800.
- 408 [10] P. Kumar *et al.*, "Integrated genomic view of SARS-CoV-2 in India," *Wellcome Open Res.*, vol. 5, p.
409 184, 2020, doi: 10.12688/wellcomeopenres.16119.1.
- 410 [11] A. Maitra *et al.*, "Mutations in SARS-CoV-2 viral RNA identified in Eastern India: Possible
411 implications for the ongoing outbreak in India and impact on viral structure and host susceptibility," *J.*
412 *Biosci.*, vol. 45, no. 1, pp. 1–18, 2020, doi: 10.1007/s12038-020-00046-1.
- 413 [12] R. C. Bhoyar *et al.*, "High throughput detection and genetic epidemiology of SARS-CoV-2 using
414 COVIDSeq next generation sequencing," *bioRxiv*, p. 2020.08.10.242677, 2020, doi:
415 10.1101/2020.08.10.242677.
- 416 [13] "Pandemic coronavirus causing COVID-19," p. 2020, 2020.
- 417 [14] S. Banu *et al.*, "A distinct phylogenetic cluster of Indian SARS-CoV-2 isolates," *bioRxiv*, p.
418 2020.05.31.126136, 2020, doi: 10.1101/2020.05.31.126136.
- 419 [15] "ICMR_Guidelines.pdf" .
- 420 [16] V. M. Corman *et al.*, "Detection of 2019 novel coronavirus (2019-nCoV) by real-time RT-PCR,"
421 *Eurosurveillance*, vol. 25, no. 3, pp. 1–8, 2020, doi: 10.2807/1560-7917.ES.2020.25.3.2000045.
- 422 [17] K. Itokawa, T. Sekizuka, M. Hashino, R. Tanaka, and M. Kuroda, "A proposal of alternative primers for
423 the ARTIC Network's multiplex PCR to improve coverage of SARS-CoV-2 genome sequencing,"
424 *bioRxiv*, no. December 2019, p. 2020.03.10.985150, 2020, doi: 10.1101/2020.03.10.985150.
- 425 [18] N. Loman, W. Rowe, and A. Rambaut, "nCoV-2019 novel coronavirus bioinformatics protocol," pp. 1–
426 5, 2020, [Online]. Available: <https://artic.network/ncov-2019/ncov2019-bioinformatics-sop.html>.
- 427 [19] C. Wang *et al.*, "The establishment of reference sequence for SARS-CoV-2 and variation analysis," *J.*
428 *Med. Virol.*, vol. 92, no. 6, pp. 667–674, 2020, doi: 10.1002/jmv.25762.
- 429 [20] Nextstrain, "Genomic epidemiology of novel coronavirus - Global subsampling," *Nextstrain.Org*, p. 1,
430 2020.
- 431 [21] K. Katoh, K. Misawa, K. I. Kuma, and T. Miyata, "MAFFT: A novel method for rapid multiple
432 sequence alignment based on fast Fourier transform," *Nucleic Acids Res.*, vol. 30, no. 14, pp. 3059–
433 3066, 2002, doi: 10.1093/nar/gkf436.
- 434 [22] L. T. Nguyen, H. A. Schmidt, A. Von Haeseler, and B. Q. Minh, "IQ-TREE: A fast and effective
435 stochastic algorithm for estimating maximum-likelihood phylogenies," *Mol. Biol. Evol.*, vol. 32, no. 1,
436 pp. 268–274, 2015, doi: 10.1093/molbev/msu300.
- 437 [23] "Auspice : An Open-source Interactive Tool for Visualising Phylogenomic Data," pp. 7–9, 2020.

- 438 [24] A. J. Drummond, A. Rambaut, and R. R. Bouckaert, “Divergence Dating Tutorial with BEAST 2.0,”
439 *Beast 2.0*, p. 19, 2013.
- 440 [25] P. Sagulenko and R. Neher, “TreeTime Documentation,” 2020.
- 441 [26] A. J. Page *et al.*, “SNP-sites: rapid efficient extraction of SNPs from multi-FASTA alignments,”
442 *Microb. genomics*, vol. 2, no. 4, p. e000056, 2016, doi: 10.1099/mgen.0.000056.
- 443 [27] P. Cingolani *et al.*, “A program for annotating and predicting the effects of single nucleotide
444 polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3,”
445 *Fly (Austin)*, vol. 6, no. 2, pp. 80–92, 2012, doi: 10.4161/fly.19695.
- 446 [28] N. L. Sim, P. Kumar, J. Hu, S. Henikoff, G. Schneider, and P. C. Ng, “SIFT web server: Predicting
447 effects of amino acid substitutions on proteins,” *Nucleic Acids Res.*, vol. 40, no. W1, pp. 452–457, 2012,
448 doi: 10.1093/nar/gks539.
- 449 [29] R. C. Team, “The R Project for Statistical Computing,” *Http://Www.R-Project.Org/*, vol. 3, pp. 1–12,
450 2013.
- 451 [30] C. Issues, A. Projects, and W. Secu, “About Nextstrain About Augur Quickstart,” pp. 1–4, 2020.
- 452 [31] E. B. Hodcroft, J. Hadfield, R. A. Neher, and T. Bedford, “Year-letter genetic clade naming for SARS-
453 CoV-2 on nextstrain.org,” *Nextstrain*, p. 2020, 2020, [Online]. Available: [https://virological.org/t/year-](https://virological.org/t/year-letter-genetic-clade-naming-for-sars-cov-2-on-nextstain-org/498)
454 [letter-genetic-clade-naming-for-sars-cov-2-on-nextstain-org/498](https://virological.org/t/year-letter-genetic-clade-naming-for-sars-cov-2-on-nextstain-org/498).
- 455 [32] K. G. Andersen, A. Rambaut, W. I. Lipkin, E. C. Holmes, and R. F. Garry, “The proximal origin of
456 SARS-CoV-2,” *Nat. Med.*, vol. 26, no. 4, pp. 450–452, 2020, doi: 10.1038/s41591-020-0820-9.
- 457 [33] D. E. Gordon *et al.*, “A SARS-CoV-2 protein interaction map reveals targets for drug repurposing,”
458 *Nature*, vol. 583, no. July, 2020, doi: 10.1038/s41586-020-2286-9.
- 459 [34] A. Pekosz, S. R. Schaecher, M. S. Diamond, D. H. Fremont, A. C. Sims, and R. S. Baric, “Structure,
460 expression, and intracellular localization of the SARS-CoV accessory proteins 7a and 7b,” *Adv. Exp.*
461 *Med. Biol.*, vol. 581, pp. 115–120, 2006, doi: 10.1007/978-0-387-33012-9_20.
- 462 [35] Q. He, Y. Li, L. Zhou, X. Ge, X. Guo, and H. Yang, “Both Nsp1 β and Nsp11 are responsible for
463 differential TNF- α production induced by porcine reproductive and respiratory syndrome virus strains
464 with different pathogenicity in vitro,” *Virus Res.*, vol. 201, pp. 32–40, 2015, doi:
465 10.1016/j.virusres.2015.02.014.
- 466 [36] Y. Sun, D. Li, S. Giri, S. G. Prasanth, and D. Yoo, “Differential host cell gene expression and regulation
467 of cell cycle progression by nonstructural protein 11 of porcine reproductive and respiratory syndrome
468 virus,” *Biomed Res. Int.*, vol. 2014, 2014, doi: 10.1155/2014/430508.
- 469 [37] L. Mousavizadeh and S. Ghasemi, “Genotype and phenotype of COVID-19: Their roles in
470 pathogenesis,” *J. Microbiol. Immunol. Infect.*, no. xxxx, pp. 0–4, 2020, doi: 10.1016/j.jmii.2020.03.022.
- 471 [38] D. X. Liu, T. S. Fung, K. K. L. Chong, A. Shukla, and R. Hilgenfeld, “Accessory proteins of SARS-
472 CoV and other coronaviruses,” *Antiviral Res.*, vol. 109, no. 1, pp. 97–109, 2014, doi:
473 10.1016/j.antiviral.2014.06.013.
- 474 [39] M. Frieman and R. Baric, “Mechanisms of Severe Acute Respiratory Syndrome Pathogenesis and Innate
475 Immunomodulation,” *Microbiol. Mol. Biol. Rev.*, vol. 72, no. 4, pp. 672–685, 2008, doi:
476 10.1128/mnbr.00015-08.
- 477 [40] D. S. S. Hassan, P. P. Choudhury, B. Roy, and S. S. Jana, “Missense Mutations in SARS-CoV2
478 Genomes from Indian Patients,” *SSRN Electron. J.*, 2020, doi: 10.2139/ssrn.3609566.
- 479 [41] R. McBride, M. van Zyl, and B. C. Fielding, “The coronavirus nucleocapsid is a multifunctional
480 protein,” *Viruses*, vol. 6, no. 8, pp. 2991–3018, 2014, doi: 10.3390/v6082991.
- 481 [42] Y. Cong *et al.*, “Nucleocapsid Protein Recruitment to Replication-Transcription Complexes Plays a
482 Crucial Role in Coronaviral Life Cycle,” *J. Virol.*, vol. 94, no. 4, pp. 1–21, 2019, doi:

483 10.1128/jvi.01925-19.

484 [43] I. Imbert, E. J. Snijder, M. Dimitrova, J. C. Guillemot, P. Lécine, and B. Canard, “The SARS-
485 Coronavirus PLnc domain of nsp3 as a replication/transcription scaffolding protein,” *Virus Res.*, vol.
486 133, no. 2, pp. 136–148, 2008, doi: 10.1016/j.virusres.2007.11.017.

487

Fig. 1a

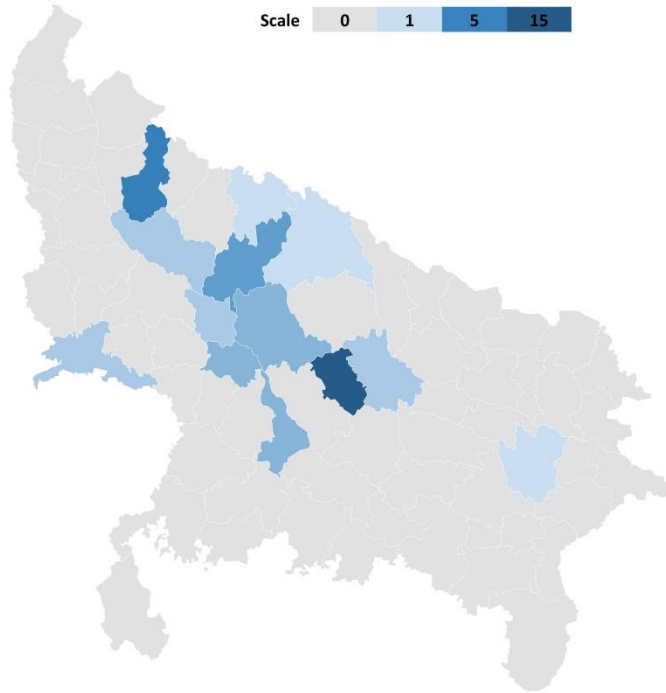
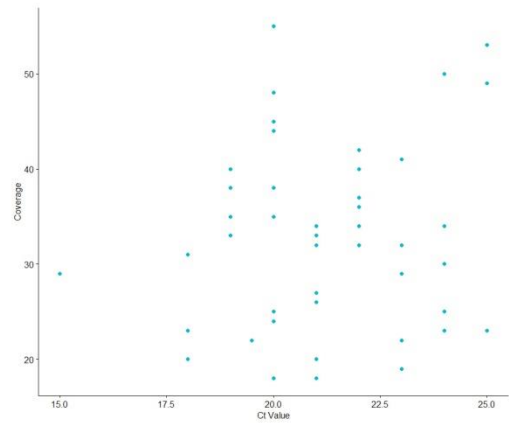


Fig. 1b



488

Fig. 2a

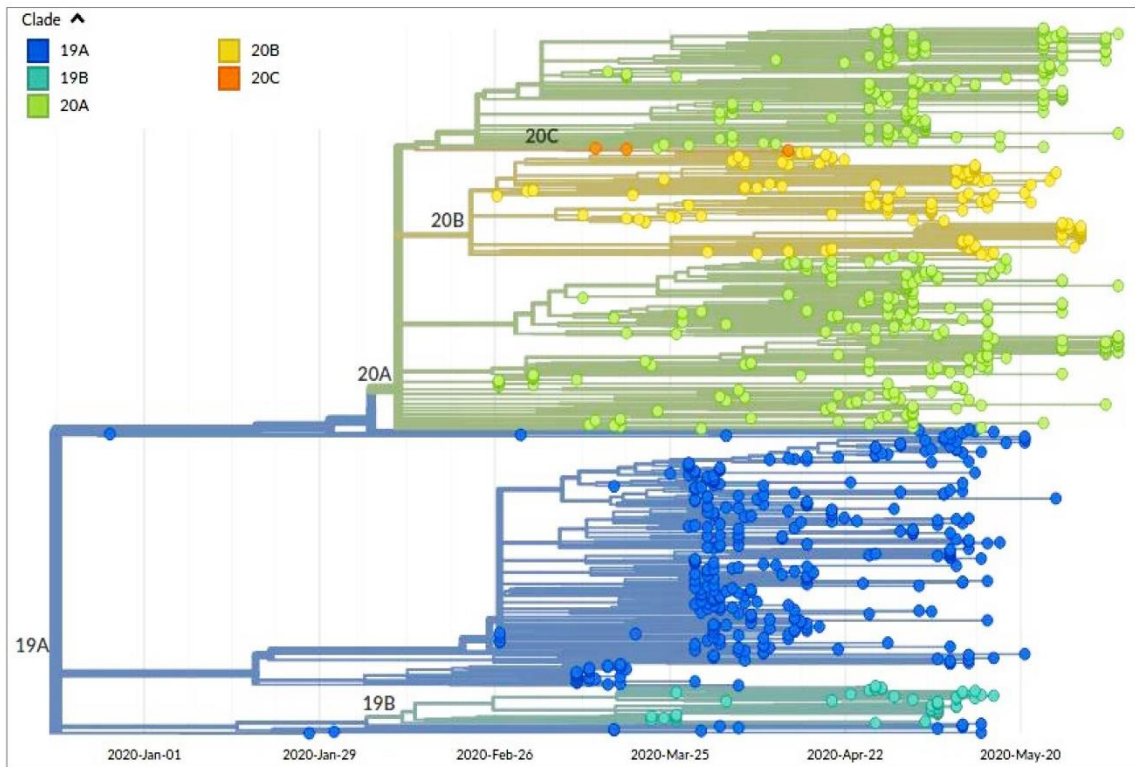
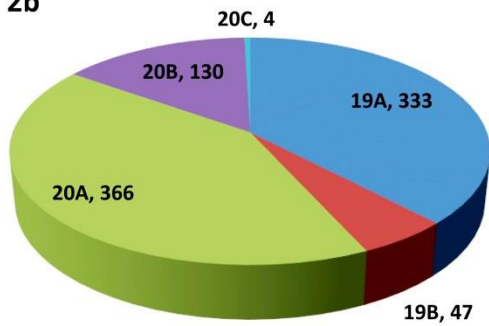


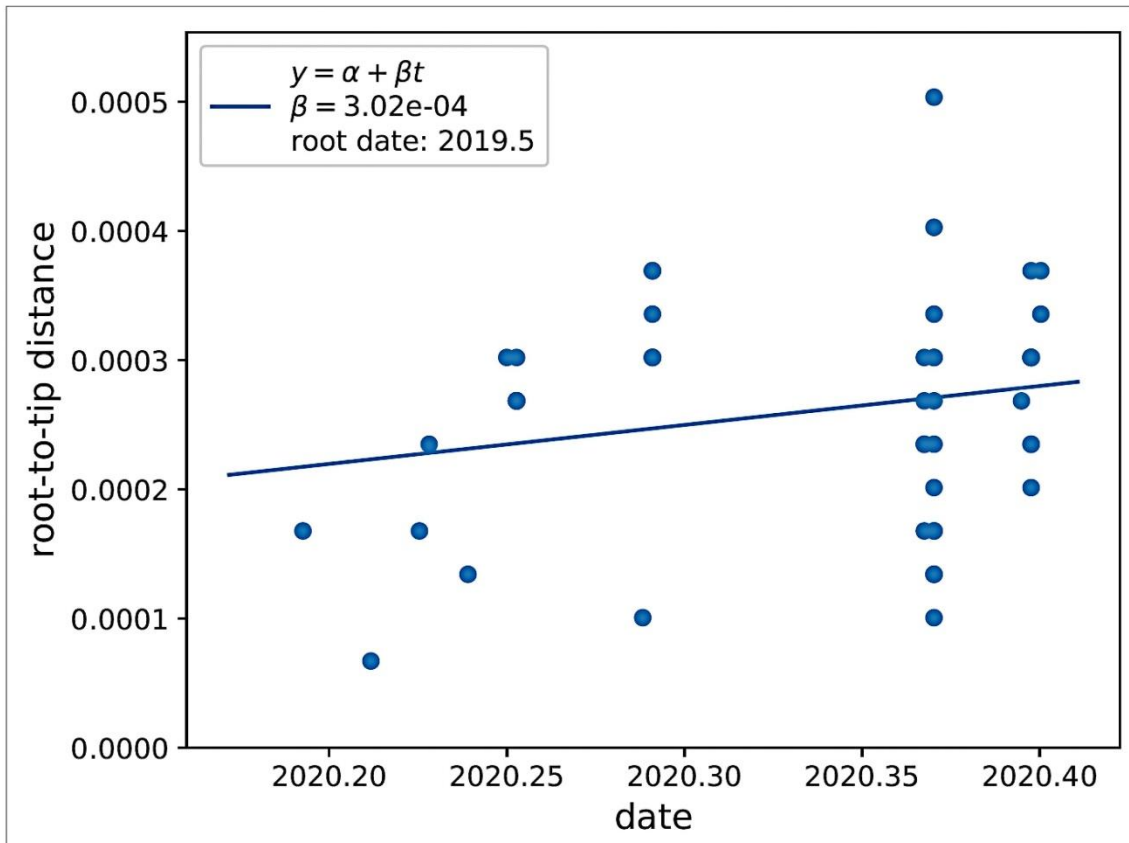
Fig. 2b



489

490

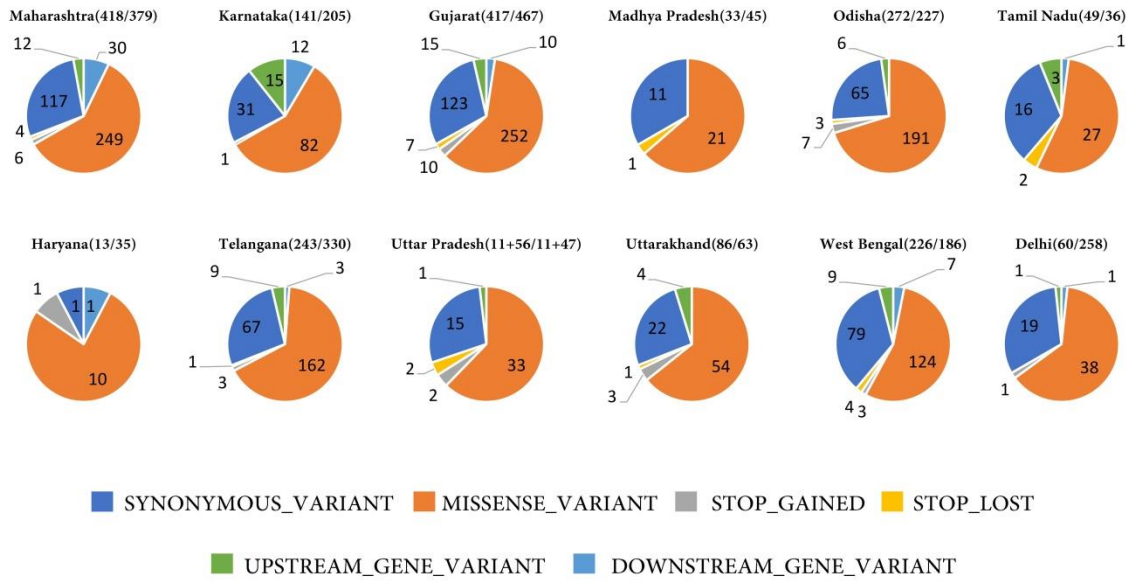
Fig. 3



491

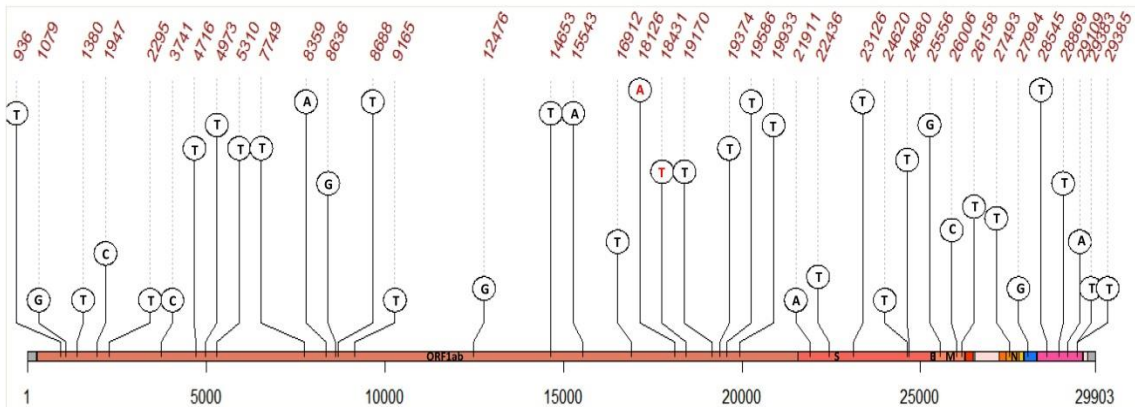
492

Fig. 4



493
494
495
496
497

Fig. 5



498
499
500
501
502

503

504 **Table -1**

State	Genome Sequenced	Unique mutation	Mutation rate	Fatality rate
Maharashtra	379	422	111.3456464	3.261071861
Gujarat	467	417	89.29336188	3.336060568
Odisha	227	273	120.2643172	0.5883102
Telangana	330	243	73.63636364	0.717308726
West Bengal	186	227	122.0430108	2.011953626
Karnataka	205	141	68.7804878	1.68566019
Uttarakhand	63	87	138.0952381	1.322401481
Delhi	258	60	23.25581395	2.663099352
Tamil Nadu	36	49	136.1111111	1.717780091
Madhya Pradesh	45	33	73.33333333	2.431820632
Haryana	35	13	37.14285714	1.10874122
Uttar Pradesh	11 + 47	11 + 56	96.55172414	1.558198114

505

506 **Table – 2**

Sample_ID	POS	REF	ALT	Gene	Mutation	Type of mutation
NBRI-N24	936	C	T	ORF1ab	Thr224Ile	missense_variant
NBRI-N25	1079	A	G	ORF1ab	Asn272Asp	missense_variant
NBRI-N38	1380	C	T	ORF1ab	Ala372Val	missense_variant
NBRI-N48	1947	T	C	ORF1ab	Val561Ala	missense_variant
NBRI-N51	2295	C	T	ORF1ab	Thr677Ile	missense_variant
NBRI-N22	3741	T	C	ORF1ab	Ile1159Thr	missense_variant
NBRI-N55	4716	C	T	ORF1ab	Thr1484Ile	missense_variant
NBRI-N42	4973	G	T	ORF1ab	Val1570Leu	missense_variant
NBRI-N47	5310	C	T	ORF1ab	Thr1682Ile	missense_variant
NBRI-15	7749	C	T	ORF1ab	Thr2495Ile	missense_variant
NBRI-N39	8359	T	A	ORF1ab	His2698Gln	missense_variant
NBRI-N24	8636	A	G	ORF1ab	Thr2791Ala	missense_variant
NBRI-N26	8688	G	T	ORF1ab	Gly2808Val	missense_variant
NBRI-N42	9165	C	T	ORF1ab	Thr2967Ile	missense_variant
NBRI-N44	12476	A	G	ORF1ab	Met4071Val	missense_variant
NBRI-N30	14653	G	T	ORF1ab	Met4796Ile	missense_variant
NBRI-N48	15543	G	A	ORF1ab	Arg5093Gln	missense_variant
NBRI-N21	16912	G	T	ORF1ab	Leu5549Phe	missense_variant

NBRI-N49	19170	C	T	ORF1ab	Ser6302Leu	missense_variant
NBRI-N50	19374	C	T	ORF1ab	Ser6370Phe	missense_variant
NBRI-N21	19586	C	T	ORF1ab	Leu6441Phe	missense_variant
NBRI-N58	21911	C	A	S	Leu117Ile	missense_variant
NBRI-N32	22436	G	T	S	Ala292Ser	missense_variant
NBRI-N49	23126	G	T	S	Ala522Ser	missense_variant
NBRI-N48	24620	G	T	S	Ala1020Ser	missense_variant
NBRI-N35	24680	G	T	S	Val1040Phe	missense_variant
NBRI-N48	25556	T	G	ORF3a	Val55Gly	missense_variant
NBRI-N39	26006	G	C	ORF3a	Ser205Thr	missense_variant
NBRI-N44	26158	G	T	ORF3a	Val256Phe	missense_variant
NBRI-N24	27493	C	T	ORF7a	Pro34Ser	missense_variant
NBRI-N34	27994	A	G	ORF8	Asp34Gly	missense_variant
NBRI-N20	28545	C	T	N	Thr91Ile	missense_variant
NBRI-N50	28869	C	T	N	Pro199Leu	missense_variant
NBRI-N57	29109	C	A	N	Pro279Gln	missense_variant
NBRI-N29	29383	G	T	N	Lys370Asn	missense_variant
NBRI-N34	29385	A	T	N	Asp371Val	missense_variant
NBRI-N20	18126	T	A	ORF1ab	Leu5954*	stop_gained
NBRI-N50	18431	C	T	ORF1ab	Gln6056*	stop_gained
NBRI-N28	19933	A	T	ORF1ab	Ter6556Cys ^{ext} ??	stop_lost

507

508 **Table -3**

Sample_ID	Gene	Mutation	SIFT Score
NBRI-N34	ORF1ab	P5496S	0.81
NBRI-N57		T6056I	0.03
NBRI-N15		T2495I	0.55
NBRI-N16		T2016K	0.59
NBRI-N21		V4691F	0.01
		V5550L	0.02
		V5894I	0.59
		T6441I	0.85
		S5483F	0.02
NBRI-N22		T5300I	0.13
NBRI-N23		N272D	0.06
NBRI-N25		P3395S	1
NBRI-N27		T6557S	0
NBRI-N28		V4797F	0
NBRI-N30		R24S	0.47
NBRI-N36		G7063S	0.7
		A372V	0.05
NBRI-N38		S558F	0.46
NBRI-N40		V1570L	0.02
NBRI-N42		T2967I	0.25
		T1682I	0.2
NBRI-N47			

NBRI-N51		T677I	0.36
NBRI-N55		T1484I	0.05
		V4980L	0
NBRI-N39		H2698Q	0.18
NBRI-N24		T224I	0.03
NBRI-N26		G2808V	0.14
NBRI-N17		D343H	0.19
NBRI-N34		D371V	0.01
NBRI-N29	N	K370N	0.02
		D371Y	0.01
NBRI-N50		P199L	0.11
NBRI-N57		P279Q	0
NBRI-N41		T362I	0.03
NBRI-N48		A1020S	0.39
NBRI-N32		A292S	0.64
NBRI-N35	S	V1040F	0.12
NBRI-N49		L5F, A522S	0.10,0.73
NBRI-N58		L117I	0.76
NBRI-N39	ORF3a	S205T	0.06
NBRI-N48		V55G	0
NBRI-N24	ORF7a	P34S	0
NBRI-N29		D34G	0
NBRI-N26	ORF8	V62L	0.54