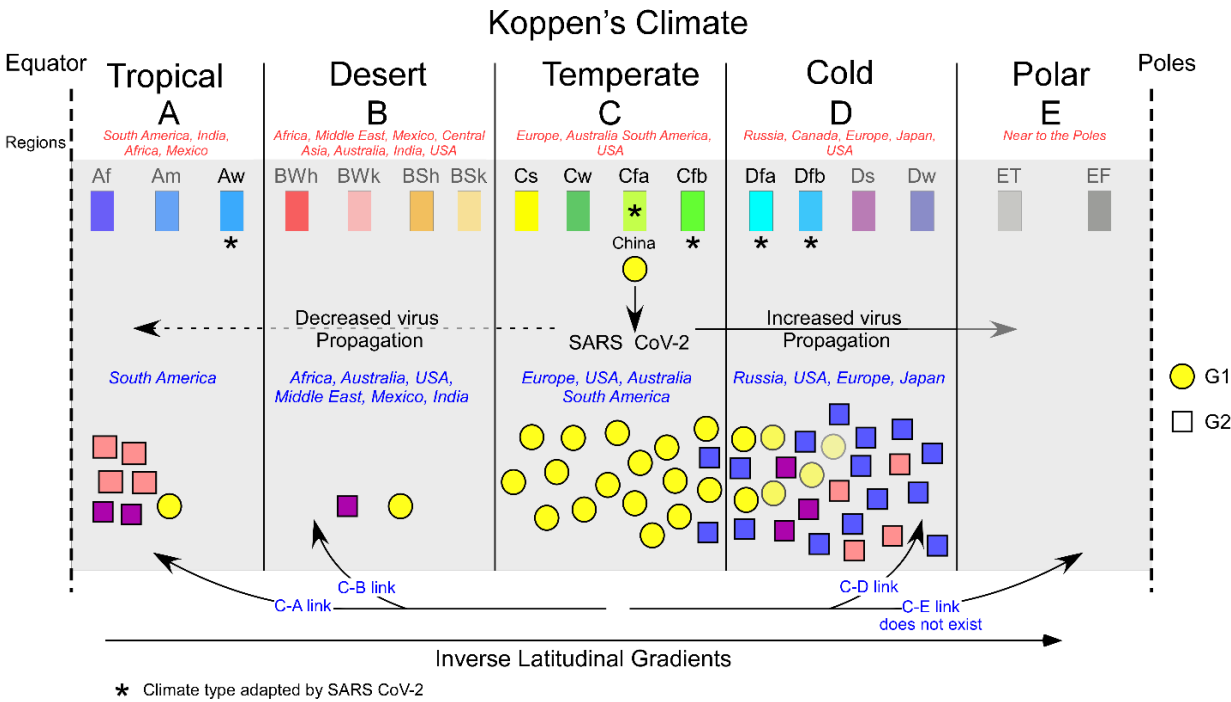# Climatic-niche evolution of SARS-CoV-2

**Authors:** Priyanka Bajaj, Prakash Chandra Arya

**Correspondence**: priyankaba@iisc.ac.in, prakasha@iisc.ac.in

## Graphical Abstract



## In Brief:

The authors elucidate adaptation of SARS-CoV-2 to different climates by studying phylogenetics & the distribution of strains on Koppen's climate map.

## Highlights:

- SARS-CoV-2 follows inverse latitudinal gradient during initial days.
- Phylogenetic network divides SARS-CoV-2 strains into two variant groups, G1 & G2.
- G1 strains is restricted to Koppen's *"temperate"* climate (mainly *Cfa-Cfb*).
- G2 strains has evolved from G1 to sustain in mainly "*humid-continental*" *(Dfa-Dfb)* & "*tropical-savannah*" *(Aw)* climate.

1

# Climatic-niche evolution of SARS-CoV-2

Priyanka Bajaj[1#*] & Prakash Chandra Arya[2#**]

1 Molecular Biophysics Unit, Indian Institute of Science, Bangalore-560012, India

2 Centre for Earth Sciences, Indian Institute of Science, Bangalore-560012, India

# Both authors have equally contributed

**Correspondence: prakasha@iisc.ac.in

*Correspondence: priyankaba@iisc.ac.in

## Abstract

COVID-19 pandemic is studied by several field experts. However, it is still unclear why it was restricted to higher latitudes during the initial days & later cascaded in the tropics. Here, we analyzed 176 SARS-CoV-2 genomes across different latitudes & climate (Koppen's climate) that provided insights about within species virus evolution & its relation to abiotic factors. Two genetically variant groups, named as G1 & G2 were identified, well defined by four mutations. The G1 group (ancestor), is mainly restricted to warm & moist, temperate climate (Koppen's C climate) while its descendent G2 group surpasses the climatic restrictions of G1, initially cascading into neighboring cold climate (D) of higher latitudes & later into hot climate of the tropics (A). It appears that the gradation of temperate climate (Cfa-Cfb) to "cold climate" (Dfa-Dfb) climate drives the evolution of G1 into G2 variant group which later adapted to tropical climate (A) as well. It seems this virus follows inverse latitudinal gradient in the beginning due to its preference towards temperate (C) & cold climate (D). Nevertheless, due to the uncertainty of COVID-19 data, the results must be cautiously interpreted & should not be extrapolated to climate types and climatic conditions other than those analyzed here for the early evolution period. Our work elucidates virus evolutionary studies combined with

45  climatic studies can provide crucial information about the pathogenesis & natural

46  spreading pathways in such outbreaks which is hard to achieve through individual

47  studies.

48  **Keywords:** SARS-CoV-2, molecular phylogeny, virus cluster SNPs, inverse

49  latitudinal gradient, climate zones, Koppen's climate.

## Introduction

51  The first case of COVID-19 pandemic caused by SARS-CoV-2 pathogen was first

52  reported from Wuhan China[1]. In spite of various precautions such as lockdown,

53  social distancing, wearing mask & sanitization, the disease was able to reach to

54  almost every part of the world infecting nearly 6 million people worldwide & putting

55  an end to nearly 370000 lives[2]. This zoonotic virus is like SARS coronavirus (79%

56  similarity) & MERS (50% similarity) & is closely related to bat derived

57  coronaviruses[1]. The SARS-CoV-2 can survive up to 3, 4, & 24 hours on aerosols,

58  copper, & cardboard respectively & up to 2-3 days on stainless-steel or plastic[1].

59  These results provide vital information about the survival of the virus in its external

60  environment, few surfaces tend to be relatively favorable than others. It has also

61  been observed that SARS-CoV-2 transmits faster than its two ancestors SARS-

62  CoV & MERS-CoV[1,3]. It is well understood that the SARS-CoV-2 has concurred a

63  larger geographical region & hosted a larger population. Since the social behaviour

64  & travelling of humans have not changed much, what makes few respiratory

65  viruses confined locally & others spread globally is still unclear. Due to its unique

66  pattern of spread, the outbreak led to a big discussion, that does climate have a

67  role in the spread of the disease. The ancestor SARS-CoV-1 losses its viability at

68  higher temperature (38°C) & relatively higher humidity (>95%)[4]. Experiments

69  support that the virus is highly stable at 4°C but is sensitive to heat[5]. The effect of

3

70   climate on COVID-19 transmission has been discussed by several authors. Studies

71   both in favour & against have been published & the topic is still debatable[6,7]. A

72   recent review compiles 61 studies relating COVID-19 with climate[1]. Several

73   climatic factors such as humidity, precipitation, radiation, temperature, & wind

74   speed affecting this virus spread have been incorporated. Both positive & negative

75   association of COVID-19 with temperature & humidity have been published[1].

76   Recently, Carlson et al. mentioned that COVID-19 transmission could be affected

77   by climate but discouraged the use of SDMs for COVID-19 transmission due to

78   their limitations, which generally takes only climatic parameters as input, as they

79   may not be appropriate tools[7]. This study was challenged by Araujo et al.

80   mentioning strengths & limitations of the tools & reasoned that $R_0$ of COVID-19

81   depends on several factors it may also be affected by climate[6]. Since only climatic

82   parameters are insufficient to capture climatic signatures of COVID-19 spread,

83   such patterns can be recognized by combining phylogenetic & climatic studies.

84   Such approach enables to probe the similarities & differences in virus genome

85   across similar & different climate types present all over the world. To understand

86   such a behaviour we have attempted to study the genomic sequence of the SARS-

87   CoV-2 across different latitudes & climate (Koppen's climate).

88   The plant & animal diversity generally decrease from equator to pole[8]. This pattern

89   is known as the latitudinal biodiversity gradient, identified & discussed by several

90   authors[8,9] with few exceptions[9]. Unlike free living plants & animals, pathogens are

91   poorly mapped & very little is known about their underlying ecological &

92   evolutionary causes[10]. Nucleotide substitution has been proposed to be one of the

93   most important mechanisms of viral evolution in nature[11]. However, factors

94   responsible for the generation of these mutations are not well understood. One of

4

95   the possible factors is adaptation to new environments dictated by natural selection

96   that discriminates among genetic variations & favours survival of the fittest[12]. Virus

97   evolution as a consequence of climate change is poorly understood. SARS-CoV-2

98   consists of large single-stranded ~30kb long positive-sense RNA. These viruses

99   majorly have a conserved genomic organization, consisting of a unique 265bp long

100  leader sequence, ORF1ab polyprotein, & structural proteins like S (spike

101  glycoprotein), E (Envelope), M (Membrane), & N (Nucleocapsid). ORF1ab

102  encodes replicase, transcriptase & helicase, essential enzymes required for

103  replication, along with non-structural & accessory proteins. Expression of non-

104  structural proteins is facilitated by ribosomal frameshifting[13]. All coronaviruses

105  express structural proteins S, E, M, N; spike glycoprotein being the most

106  immunogenic to T-cell response[14]. Spike glycoprotein of coronaviruses binds to

107  human angiotensin-converting enzyme 2 (hACE2) receptor for viral fusion & entry

108  & is the main target for neutralizing antibodies & development of vaccines[15].

109  Membrane protein is also antigenic as it stimulates a humoral immune response[16].

110  E protein is responsible for virus assembly & release of virion particles[17].

111  Nucleocapsid protein packages RNA genome into a helical ribonucleocapsid

112  protein (RNP) complex during virion assembly & is capable of eliciting an immune

113  response[18]. Since it is still not clear whether SARS-CoV-2 evolution & spread have

114  relation with climate, our study may act as a missing link between genomic

115  sequence, climate & COVID-19 severity. If SARS-CoV-2 is responding towards

116  external climate it can be delineated by superimposing its genomic variants across

117  different latitudes & Koppen's climate. The earliest & the most simple classification

118  of Earth's climate is based on latitudes which divide the Earth's climate into seven

119  climate zones, North Frigid Zone (NFZ), North Temperate Zone (NTZ), North

120    Subtropical Zone (NSTZ), Tropical Zone (TZ), South Subtropical Zone (SSTZ),

121    South Temperate Zone (STZ) & South Frigid Zone (SFZ)[19]. Wladimir Koppen

122    presented a modified classification of Earth's climate based on the precipitation &

123    temperature[20]. He divided Earth's climate into five major climates, A (Tropical), B

124    (Arid), C (Temperate), D (Cold or Continental) & E (Polar) which are further

125    subdivided into 30 climate types[20]. To understand the effect of climate on SARS-

126    CoV-2 evolution, the present study comprises of three parts, (1) latitudinal

127    distribution of COVID-19 cases, (2) sequence analysis of SARS-CoV-2 strains, (3)

128    mapping SARS-CoV-2 strains across different climates. These combined studies

129    can provide insights on within species evolution & preferential distribution of SARS-

130    CoV-2 across different climatic zones which might be difficult to probe through

131    individual studies. These results will provide useful information to design

132    efficacious vaccines which can be stored and transported in a wide range of

133    temperature and humid conditions, thereby minimizing cold storage costs.

## Results

### Distribution of COVID-19 cases across latitudes

136    For an overview of the latitudinal preference of SARS-CoV-2, we have plotted per-

137    million active cases of SARS-CoV-2 across different climate zones (Figure1a).

138    Results show that 81% of the cases belong to NTZ (30°N-66.5°N), 4% to NSTZ

139    (23.5°N-30°N), 14% lie in the TZ (23.5°N-23.5°S), 1% in the STZ (30°S-66.6°S) &

140    negligible number of cases (<0.5%) are from remaining climate zones. Statistical

141    difference exists between number of COVID cases in Temperate Zone versus

142    other climate zones (paired t-test two-tail, $P$<.001). The spread of COVID-19 is

143    dominant in the higher latitudes which is usually uncommon as majority of

6

144 terrestrial taxa prefers to stay near tropical region, suggesting that SARS-CoV-2

145 follows inverse latitudinal gradient in early stages of pandemic. Since a majority of

146 the cases lie in the NTZ, we have further divided this zone into an interval of 7°

147 latitude i.e. 30°N-37°N, 37°N-44°N, 44°N-51°N, 51°N-58°N & 58°N-66.5°N. We

148 found 9% of the cases fall in latitude range 30°N-37°N, 46% in 37°N-44°N, 21% in

149 44°N-51°N, 14% in 51°N-58°N & 10% in 58°N-66.5°N (Figure1b). The results show

150 a peak of COVID-19 cases in between 37°N to 51°N latitudes, the dominant

151 Koppen's climate between these latitudes is temperate (C) & continental climate

152 (D). The general characteristics of these climates are prevalence of high

153 atmospheric circulation with anticyclones during winters, with an average

154 temperature of ~15°C for C & ~< 10°C for D climate, with relative humidity ranging

155 between ~50-80%. Since the major distribution of SARS-CoV-2 is confined within

156 a latitude range, this trend could be random or there might be a strong underlying

157 cause driven by underlying principles. The latitudes have a very high control on

158 climate, a detailed investigation of the Koppen's climate under the light of genomic

159 sequences is carried out to understand the distribution pattern across the globe.

160 **Molecular phylogeny analysis to infer genomic similarities & their**

161 **distribution in different climates**

162 To probe genomic similarities between SARS-CoV-2 virus isolates, a phylogenetic

163 tree was constructed by aligning 176 virus genomes to the reference genome[21]

164 retrieved from GISAID. Interestingly, our Multiple Sequence Alignment (MSA)

165 results reveal sixty virus cluster Single Nucleotide Polymorphisms (SNPs) (see

166 methodology). Table1 comprises of SNPs of virus clusters across different climatic

167 zones, Koppen's climate & climate type. Climatic parameters (temperature &

168 precipitation) for each virus strain is mentioned in TableS2. Based on phylogenetic

7

169   clustering, 176 SARS-CoV-2 strains are majorly divided into two groups, we named

170   them as G1 (1-58) & G2 (59-176) (Figure2). Predominantly four mutations

171   distinguish G2 from G1 group, i.e., 1) a synonymous mutation (C241T) appeared

172   in the unique leader sequence, 2) F924 (C3037T) appeared in nsp3, encoding for

173   papain-like proteinase[22], 3) a non-synonymous mutation, P214L (C14408T) arose

174   in ORF1b, that codes for four putative non-structural proteins (nsp13, nsp14, nsp15

175   & nsp16), functionally involved in replication-transcription complex[23], & 4) D614G

176   (A23403G) arose in S gene, encoding spike glycoprotein[14] (Figure3a). Among four

177   mutations in G2, the D614G mutation, lying in spike glycoprotein was widely

178   studied due to its higher infectivity & involvement in entering the host cell through

179   hACE2 receptors[24–27]. The other three mutations in G2 have co-evolved with

180   D614G making it distinguishable from G1. We explored the extent of genome-wide

181   divergence of G1 & G2 group across different climate zones & Koppen's climate

182   (Figure3b). 59% of G1 viruses fall in NTZ, 14% in NSTZ, 12% in TZ, 10% in SSTZ

183   & 5% in STZ. 76% of the virus isolates in G2 group are present in the NTZ, 13.5%

184   in TZ, 7.6% in STZ & remaining 2% is equally distributed in NSTZ & SSTZ, showing

185   G2 strain variants evolved to adapt to temperate zones as their population

186   decreased drastically in the subtropical zones. These results show both G1 & G2

187   strains have a strong preference towards higher latitudes i.e., NTZ, which agrees

188   with the analyzed worldometer data (Figure3c). It also supports that in the initial

189   stages of the pandemic, the virus isolates follow inverse latitudinal gradient.

190   Mapping viral strains on Koppen's map (thoroughly discussed in the next section)

191   reveal their prevalence majorly in the C & D climate (Figure3d). 71% of G1 lie in C

192   climate, 17% in D & the remaining is equally distributed in the A & B climate. 54%

193   of G2 lie in C climate, 36% in D, 9% in A & 1% in B climate pointing towards a

preferential shift of the novel coronavirus towards D climate (Figure3b), alluding G2 is climatically & genomically more diverse than G1. The analysis suggests that the G1 group is mostly restricted to temperate climate (C) & G2 is climatically & geography widely distributed, it is possible that these mutations were acquired by G1 to stabilize itself in different climates hence allowing it to spread globally. Similar climatic concordance with the temperate climate (C) was also observed for SARS-CoV that was responsible for 2002-2004 epidemic as it prevailed in regions of Australia, Europe, Canada and, China[28], having Koppen's C climate. Such similar occurrence of SARS-CoV and G1 group of SARS-CoV-2 hints towards, why initially G1 variant group (consisting of the reference genome NC_045512[21]) that has 79% similarity to SARS-CoV[1] was majorly located in the temperate climate and latter it evolved into G2 variant group that allowed it to extend its climatic boundaries into temperate, cold and tropical climate. These results suggest these four SNPs could be the key factors in increasing the virulence, transmission & sustainability of the virus in humans.

We further analyzed the order in which the phylogenetic clusters evolved from the ancestor 45-57 cluster (containing the reference genome, Strain ID: 50) based on nodes, mutational branches & branch length. The order in which the virus evolved is 44-47 (G1440A, G392D; G2891A, A876T), 1-22 (C8782T, S2839; T28144C, L84S), 33-43 (G26144T, G251V), 23-32, 58-61 (C15324T, N519), 80-115 (G28881A, G28882A, R203K; G28883C, G204R), 116-125 (A20268G, L2167), 126-176 (G25563T, Q57H) & 62-79 (cluster, acquired genomic mutation & its corresponding amino acid mutation). In Figure3e, looking at the distribution of the viruses in different climate zones, no such preference was observed as the virus evolved. Virus cluster 58-61, linking G1 & G2 has an equal distribution of virus

219    strains in C & D climate. The virus cluster 80-115 of G2 more closely related to G1,

220    is widely distributed in A, C & D climates. Within 80-115 virus cluster, 106-115

221    cluster shows distribution in C & A climate. A trend was observed that virus clusters

222    in G2 group gradually evolved to sustain in Koppen's D climate which supports our

223    previous observation. These analyses led us interpret G1 group (ancestor), is

224    mainly restricted to warm & moist, temperate climate (C) while its descendent G2

225    group surpasses the climatic restrictions of G1, initially cascading into neighboring

226    cold climate (D) of higher latitudes & later into hot climate of the tropics (A). Within

227    these major virus clusters, small clusters also exist as shown in Table1 with their

228    mutations along with their climatic distribution.

229    We have examined whether climatic conditions exhibit any selective pressure on

230    each gene (Figure3f). Since, the present picture of the data appears that SARS-

231    CoV-2 is following inverse latitudinal gradient, as expected all genes are having

232    mutations in NTZ, suggesting the virus is probably using varied mechanisms to

233    adapt to the two main climates of NTZ i.e. temperate (C) & cold or continental (D).

234    Mutations in the M gene are only pertaining to NTZ & NSTZ & are present in C &

235    D climate. In particular, there is a surge in the virus strains carrying SNPs in ORF8

236    in the NSTZ (20%). 77% of the SNPs in ORF8 lie in the C & 20% in the D climate.

237    Overall, the distribution of virus cluster SNPs of ORF1ab, S, ORF3a, & N gene

238    follows a similar pattern across all the climatic zones & Koppen's climate, implying

239    no difference in selective pressure of the climate in generating mutations in these

240    genes. S, M, & N proteins are immunogenic[14,16,18], implicating virus evades

241    immune response by introducing these substitutions.

242    Apart from non-synonymous mutations, synonymous mutations within the gene

243    can also significantly affect protein function due to codon usage bias[29] & through

10

244    mechanisms such as ribosome stalling[30] & mRNA secondary structure formation[31].

245    We probed the frequency of derived synonymous versus non-synonymous

246    mutations & observed a very similar distribution pattern of the derived synonymous

247    versus missense mutations across all climate zones & Koppen's climate

248    (Figure3g). These analyses suggest novel coronavirus is using varied mechanisms

249    both at the transcriptional as well as translational level to adapt, survive, & increase

250    infectivity in all types of climates. These findings unequivocally bolster a

251    requirement for further prompt, comprehensive studies that join genomic

252    information, epidemiological information, & climatic distribution with COVID-19

253    severity.

254    **Distribution of strains across Koppen's climate**

255    To probe the relation between climate & SARS-CoV-2 strains, we superimposed

256    genomic information along with their geolocations on the climate map of Wladimir

257    Koppen (Figure4). We carefully examined the distribution of strains on Koppen's

258    map & an overview of the map shows, the distribution of 176 strains are mainly

259    concentrated in the western coasts of Europe & North America, & eastern coasts

260    of China, North America, Australia & South America (Figure4). Throughout the text

261    Koppen's climate type is marked within quotations & its standard symbol is written

262    within brackets e.g., "*humid-subtropical*" *(Cfa)*.   List of Koppen's symbol of each

263    climate type is given in Supplementary TableS3 & its criteria for classification is

264    given in Supplementary TableS4. Mostly the SARS-CoV-2 strains are distributed

265    in the "*humid-subtropical*" *(Cfa)* & "*marine-temperate*" *(Cfb)* & "*humid-continental*"

266    *(Dfa-Dfb)* climate &, two strains from virus clusters (80-115 & 126-176) belonging

267    to South America, are found in "tropical-savanna" *(Aw)* of 'A' climate

268    (Supplementary TableS5). The map displays ~86% (n=176) of virus isolates are

269    distributed in the coastal regions & the remaining in continental region (Chi-square

270    test, *P*<.001, Figure5a). Around ~73.86% of the total strains are distributed in

271    "humid-subtropical" (Cfa) & "marine-temperate" (Cfb) climate type of C climate &

272    "humid-continental" (Dfa-Dfb) climate type of D climate. The remaining ~26.14%

273    strains are distributed in other climate types of other Koppen's climate including

274    non 'Cfa-Cfb' of C climate & non 'Dfa-Dfb' of D climate (Figure5b). It seems that

275    spread of COVID-19 is maximally in areas with 'Cfa' & 'Cfb' climate type.  The

276    climatic parameters (temperature & precipitation) in which these strains were found

277    were analyzed. Statistically, significant difference was found in the latitudes of G1

278    ~24.14±3.5 (mean±s.e.) & G2   ~34.03±2.7 (one-way ANOVA, *P*=.03251,

279    Figure6a). Statistically, significant difference was observed in the temperatures of

280    G1 (15.82±0.75 °C (mean±s.e.) & G2 (11.67±0.68 °C) strains (one-way ANOVA,

281    *P*<.001, Figure6b). However, the difference in precipitation for G1 (1046.95±80

282    mm) & G2 (896.64±35.48 mm) strains is statistically not significant (one-way

283    ANOVA, *P*=.06118, Figure6c). The latitudes & temperature are inversely related to

284    each other (r = -0.6649, Supplementary Figure1a & Figure6d), which explains the

285    occurrence of G1 strains in lower & G2 strains in higher latitudes. Such relation

286    between latitude & precipitation has not been observed (r = -0.3064,

287    Supplementary FigureS1b & Figure6e). A mesh plot simultaneously evaluates all

288    climatic parameters for both G1 & G2 strains, the results agrees to the limited

289    temperature & wider precipitation range of G1 group & interestingly the G2 group

290    appears in a wider temperature & shows a preferential shift towards lower

291    temperature which is evident from the fact that it initially appeared more in higher

292    latitudes (Supplementary FigureS2). A complete description of the distribution of

293    G1 & G2 strains lying in different countries &/or continents of the world is provided

294    in Supplementary FigureS3 & Supplementary Material.

## Discussion

296    Unlike majority of terrestrial organisms, it appears that SARS-CoV-2 is following

297    an inverse latitudinal gradient, a possible explanation from our analysis is that

298    evolution from G1 to G2 might help to sustain this virus in temperate & cold climates

299    which might change with time as it appears it is adapting to different climate. Mainly

300    four mutations in leader sequence, ORF1ab, & S gene were identified that led G1

301    to evolve into G2. The leader sequence & ORF1ab is involved in replication &

302    transcription, & the S gene is involved in binding to the host cell through hACE2

303    receptors. Substitutions in the ORF1ab gene may increase the synthesis of

304    replicase-transcriptase complex, thus, increasing the replication rate of the virus &

305    blocking the host innate-immune response. 614 position in spike glycoprotein lies

306    near the S1/S2 subunit junction where the furin-cleavage site is present (R667)

307    that enhances virion cell-cell fusion[32]. This suggests, aspartate to glycine

308    substitution in the vicinity of the furin-recognition site may result in a conformational

309    change of the spike glycoprotein that favors higher affinity of the Receptor Binding

310    Domain (RBD) to hACE2. A recent article showed retroviruses pseudotyped with

311    Glycine at 614 position infected ACE2-expressing cells markedly more efficiently

312    than those with Aspartic acid due to less S1 shedding & greater incorporation of

313    the S protein into the pseudovirion[24]. Several studies reported D614G mutation is

314    increasing at an alarming rate[25,26]. Few observed that this alteration correlated with

315    increased viral loads in COVID-19 patients[25]. This is consistent with the

316    epidemiological data showing proportion of viruses bearing G614 is correlated to

317   increased case fatality rate on a country by country basis[27]. This substitution

318   coevolved with substitution in the leader sequence, nsp3 & RdRp proteins,

319   suggesting these mutations allow the virus to transmit more efficiently. This

320   explains these mutations have not emerged merely because of founder's effect but

321   this virus under selection pressure has made itself more stable & infective. Also,

322   Forster et al. observed in his phylogenetic analysis the preferential geographical

323   spread of SARS-CoV-2 & provided a plausible cause which could be founders

324   effect or immunological or environmental effect[33]. Although there is a possibility

325   that the stable variant might have appeared because of host innate immune

326   response or some unknown reason, in such a case it would not show any close

327   association with climate. Since our analysis shows largely G1 is restricted to

328   temperate climate & G2 spreads in temperate & adjoining climates, a pattern that

329   is consistently observed all over the world, led us doubt that climate could be one

330   of the possible selective pressures towards which SARS-CoV-2 responds by

331   altering its genome. Through our analysis we are inclined to say that climate does

332   not display any selective pressure on each gene of SARS-CoV-2.  Our genomic

333   analysis of virus strains show that the novel coronavirus undergoes both

334   synonymous as well as non-synonymous mutations throughout its genome in

335   various climates, suggesting the novel coronavirus uses multiple mechanisms both

336   at the transcriptional & translational level for evading the immune response,

337   developing drug resistance & increasing pathogenesis. However, the actual role of

338   these mutations is not yet determined, & these studies need to be further

339   enlightened by biophysical & biochemical studies. Such mutational insights will aid

340   the design of efficacious vaccines that can be stored and transported in a wide

341   range of temperature and conditions, thereby minimizing cold storage costs.

342    To delineate the signatures of underlying abiotic factors (temperature,

343    precipitation, & latitude) responsible for evolution of SARS-CoV-2 (n=176),

344    spreading patterns of G1 & G2 strains were carefully examined on Koppen's

345    climate map. Figure4 shows an elevated spread of COVID-19 in the western &

346    eastern coasts of the continents & a diminished spread in the hot & cold deserts.

347    The G1 strains are majorly present in the eastern & western coasts of the

348    continents & G2 strains lie in both the coastal regions & continent's interior. On a

349    closer inspection, the eastern coasts of continents consist of "humid-subtropical"

350    (Cfa) climate while the western coasts of continents consist of "marine-temperate"

351    (Cfb) commonly known as east & west coast climate, respectively. These two

352    climates are very similar to each other & belongs to temperate climate also known

353    as C type climate of Koppen's classification scheme. A very large portion (~94%)

354    of habitable China consists of temperate climate (C), i.e., humid subtropical climate

355    (Cfa), which explains presence of only G1 strains in China & one strain of G1 is

356    present in cold climate (D) present near the transition of temperate (C) to cold

357    climate (D), thus probably temperate climate was suitable for G1. A similar

358    association of G1 with temperate climate (C) was found in eastern & western coast

359    of North America, eastern coast of South America, western coast of Europe &

360    eastern & western coast of Australia. Statistically, distribution of G1 strains all over

361    the globe is in concordance with the temperate climate & strongly favor C climate

362    (Chi-square test, $P$<.001) as compared to any other climate. If climate does not

363    have any role in the evolution & preferential spread of coronavirus, in such a case

364    G1 would have been evenly distributed in all climate types which is not the case.

365    Few exceptions of G1 seen in other climate types are most probably because of

366    travel as they remained subsided in that climate, implying their inability to sustain

15

367    in other climate types. Such a significant difference in association of climate with

368    G1 & G2 population could not merely arise due to human transportation. It appears

369    that the G1 strains existed in temperate climate all over the world but could not

370    extend their geographical territories beyond temperate climate. Contrastingly, the

371    evolved G2 strains can sustain in temperate (C), cold (D) & tropical (A) climate. It

372    appears that G2 strains enters the continent's interior through D climate (e.g., North

373    America & Russia). Temperate climate (C) generally grades into cold climate (D)

374    & deserts (B) in the northern hemisphere (e.g., C to D: Europe to Russia, & USA

375    to Canada; C to B: China, & USA). In southern hemisphere, gradation of temperate

376    climate (C) into tropical climate (A) & deserts (B) exists (e.g., C to A, Brazil; C to B,

377    Australia), C to A transition is identified by virus cluster 105-115 in phylogenetic

378    tree. In Russia, 91.3% (21/23) of the strains belong to G2 (Figure4), are mainly

379    present in the ~8500 km long & 600-1700 km wide D climate belt ('Dfa-Dfb-Dw'),

380    suggesting the G2 strains might have adapted to the D climate (Chi-square test,

381    $P$<.001). Similar observations are seen for North America, South America &

382    Australia. The eastern & western coasts of North America have temperate climate

383    & are connected by cold climate along USA-Canada boundary (i.e., having humid

384    subtropical (Cfa) in eastern coast & marine temperate climate (Cfb) in western

385    coast) (Figure4). The G2 strains follows this cold climate (Dfa-Dfb) belt which is

386    ~3800 km long & ~600 to 1000 km wide. The dominance of G2 & nearly absence

387    of G1 population in cold climate of North America is similar to the observations of

388    Russia. Our analysis suggests that a fall of temperature from temperate to cold

389    climate might have dictated the evolution of G1 into G2 variant group (Chi-square

390    test, $P$<.001). Similarly, a change in climate from C to A probably made the strains

391    stable in tropical regions.  Throughout the world, the G1 strains have expanded

392  only in temperate climate, suggesting a close relationship between G1 with

393  temperate climate, however the evolved G2 strains were able to infect temperate,

394  cold, & tropical climate. There might be other factors controlling virus evolution, but

395  the possibility of virus evolving to sustain in different climates of the Earth cannot

396  be neglected as several other studies also verified that viruses do respond towards

397  temperature[4,5]. Our analysis also shows that it is highly possible SARS-CoV-2 has

398  evolved to sustain in different climate. Studies combining genetic information with

399  climate can provides useful information about virus evolution & possible climatic

400  pathways during an outbreak.

## Conclusion

402  It is reasonable to assume COVID-19 transmission pathway & evolution is

403  influenced by climate. Phylogenetic network classified 176 SARS-CoV-2 strains

404  into two variant groups G1 & G2. The G1 strains were habituated to C climate that

405  evolved into G2 by undergoing significant mutations (C241T in leader sequence,

406  F924 in ORF1a, P214L in ORF1b & D614G in S gene), plausibly extended its

407  climatic boundaries from C to D climate, displaying role of natural selection on virus

408  evolution. In our analysis SARS-CoV-2, were found resistive to desert climate (B).

409  Gradually, strains are adapting to A climate in South America. The strains adapted

410  to *"tropical-savannah" (Aw)* climate are a threat to all the tropical countries, which

411  were initially less affected by COVID-19. There is a high possibility that the

412  evolutionary pathway adopted by the virus is temperate climate to cold climate for

413  higher latitude & temperate climate to tropical climate for lower latitudes.

414  Nevertheless, due to the uncertainty of COVID-19 data, the results should be

415  carefully interpreted & should not be extrapolated to climate types and climatic

17

416    conditions other than those analyzed here for the early evolution period. The study

417    agrees that viruses are sensitive to their environment & respond towards naturally

418    occurring abiotic factors such as temperature, latitude & humidity to sustain in

419    different climate of the Earth, which also provides insights about seasonal

420    variations possibly being a strong reason for the spread of other viral diseases as

421    well. Here we showed a more refined description of genes based on phylogenetics

422    & their distribution across different climate zones. This finer-grained analysis led to

423    highly relevant insights on evolutionary dynamics of poorly understood SARS-CoV-

424    2 genome & provides vital information about the direction of the spread & highlights

425    vulnerable regions of Earth. Such inter-disciplinary studies will play an imperative

426    role in designing antiviral strategies & taking pre-emptive precautionary measures

427    to combat COVID-19.

## Methodology

### Distribution of SARS-COV-2 across latitudes

430    The COVID-19 data is obtained from the 'worldometer' website, a trusted source

431    of COVID-19 database which provides global COVID-19 live statistics[34]. 'Active

432    cases per million population' for different countries were analyzed (assessed on 25

433    April 2020). To check the latitudinal preference of SARS-CoV-2, the countries of

434    the world were segregated based on their latitudes & per million COVID-19 cases

435    were plotted between the latitudes (90°N to 66.5°N), (66.5°N to 23.5°N), (23.5°N

436    to 23.5°S), (23.5°S to 66.5°S), & (66.5°S to 90°S). North Temperate Zone was

437    further divided in an interval of 7° latitude. Distribution of SARS-CoV-2 between

438    these latitudes was analyzed & compared.

439

**Molecular phylogenetic analysis**

185 full-length SARS-CoV-2 genomic sequences from countries across the globe, with genome length more than 29 kb & high coverage were obtained from Global Initiative on Sharing Avian Influenza Data (GISAID) database, accessed till 2 May 2020 & the reference genome was retrieved from GenBank24 (Table S1). To avoid bias related to the geographical area covered by a country, genomic sequence of strains isolated from different locations from each country was retrieved, depending on the availability of data. The sequences were aligned to the full reference genome[21] by using Biomanager & Seqinr packages of R (version 3.6.3). Among 185 genomes, some partial genomes were discarded. NC_045512 genome sequence was used as reference & the genomic coordinate in this study is based on this reference genome. Based on protein annotations, nucleotide level variants were converted into amino acid codon variants for alignments when its location within a gene was identified. The amino acid position numbering is according to its position within the specified gene (CDS) as annotated in reference sequence (NC_045512, NCBI)[21]. To ensure comparability, we trimmed the flanks of all sequences. The aligned sequences were used to construct a phylogenetic tree using MEGA X[35]. The evolutionary history was inferred using the Neighbor-Joining method (500 bootstrap tests)[36]. The optimal tree with the sum of branch length = 0.01116462 is shown. The tree is drawn to scale, with branch lengths in the same units as those of the evolutionary distances used to infer the phylogenetic tree. The evolutionary distances were computed using the Maximum Composite Likelihood method[37] & are in the units of the number of base substitutions per site. All ambiguous positions were removed for each sequence pair (pairwise deletion option). A total of 29408 positions were present in the final dataset. The results are

19

465  presented in the form of DNA sequencing i.e., U (uracil) is read as T (thymine). We

466  have labeled each virus strain by the GISAID Accession ID & the location from

467  which it was isolated in the format "Location|EPI ISL Accession ID", in the

468  constructed phylogenetic tree. For ease of visualization, we have marked a new

469  Strain ID (1 to 176) against each SARS-CoV-2 isolate in the phylogenetic tree

470  (Figure2). The same Strain ID is used for the climatic studies in this article. High-

471  frequency SNPs (Single Nucleotide Polymorphisms) distinguishing one virus

472  cluster from the others is referred to as "virus cluster SNPs" throughout this paper.

473  **Mapping virus strain on the Koppen's climate map**

474  The location of each SARS-CoV-2 strain is obtained from the METADATA file

475  provided in GISAID database for each viral isolate (Table S1). The coordinates of

476  the locations were taken from the official website of USGS Earth Explorer[38]. The

477  Gieger-Koppen's climate map is used for climatic studies[20]. The Koppen climate

478  type, temperature, precipitation of each strain is assessed from weatherbase[39] &

479  CLIMATE.ORG[40]. The map is georeferenced by using 'Arc-GIS 10.1'[41]. The

480  locations of all strains (n=176) were transferred to the georeferenced map[41]. On

481  the map, the G1 strains were symbolized as 'Yellow-circle', & G2 as 'Square'

482  (Figure4). Each strain in the map is labelled as per their Strain ID (1 to 176)

483  (Figure4), the map combines information of the phylogeny, climate, & global

484  distribution of SARS-CoV-2. These locations were classified into coastal &

485  continental region, we define the coastal region as land region < 500 km from the

486  ocean/sea & the continental region as land lying >500 km from the coastline

487  measured through google maps.

488

**Statistical analysis**

Two-tailed paired t-test & Chi-square test were performed in Microsoft Excel (2016) to test null hypothesis H1, H2, H3 & H4 related to latitudinal preference (H1), climatic preference (H2 & H3) & regional preference (H4) of SARS-CoV-2. H1: SARS-CoV-2 follows latitudinal biodiversity gradient. H2: Majority of G1 strains do not lie in temperate climate (C). H3: Majority of G2 strains do not fall in temperate (C) & cold (D) climate. H4: The virus isolates are equally distributed in coastal & continental region. Histograms depicting the distribution of coronavirus in coastal region, continental region, Koppen's climate & climate type were plotted using R (version 3.6.3). SigmaPlot10 was used to generate box plot, regression plot, & mesh plot to statistically compare frequency distribution of latitude, temperature, & precipitation of G1 & G2 strains. We performed one-way ANOVA to estimate statistical differences in the latitude, temperature & precipitation between G1 & G2 virus populations. Various scatterplots between latitude, temperature, & precipitation of G1 & G2 strains were plotted in R (version 3.6.3). Values were considered statistically significant for P values below 0.05. Exact P values are provided in appropriate figures.

**Data accessibility**

The full-length genomic sequences were downloaded from GIS-AID website (https://www.gisaid.org/), an open source database for influenza viruses. The data is downloaded as FASTA file along with the acknowledgement. The location of each strain is accessed from its METADATA file. The Koppen's Climate map is taken from the published article[20] (Peel et al., 2007). The Koppen climate type, temperature & precipitation for each strain is taken from weatherbase

513    (https://www.weatherbase.com/) & CLIMATEDATA.ORG (https://en.climate-

514    data.org/). The total cases of COVID-19 all over the globe is retrieved from

515    worldometer website (https://www.worldometers.info/coronavirus/). Refer

516    Supplementary Tables S1-S5. The code is available from the corresponding author

517    on request.

518    **Potential caveats**

519    We acknowledge several caveats about our analyses. Our data from the tropics is

520    limited because at the time of data collection (SARS-CoV-2 strains) from all over

521    the world, the strains from the tropical countries were very limited, from few tropical

522    regions strains were available (e.g., Ghana (Africa); India, Mexico, Nepal,

523    Pakistan) but the data has been discarded due to the travel history of the strains,

524    a large fraction of strains without travel history have large gaps in genomic

525    sequences which were not suitable for the present study. Also, case history of each

526    patient is not reported in the METADATA file as collecting all information from each

527    patient is time-consuming. Hence, there are chances patients from whom these

528    strains were isolated may have a migratory history. Data from different individual

529    locations without travel history & large gaps in genomic sequences have been

530    incorporated in this study. To overcome this, the inverse latitude gradients were

531    studied based on the total number of COVID-19 cases all over the globe. Due to

532    the uncertainty of COVID-19 data, these results should be carefully interpreted &

533    should not be extrapolated to climate types and climatic conditions other than those

534    analyzed here for the early evolution period.

535

536

# References

1.  Briz-Redón, Á. & Serrano-Aroca, Á. The effect of climate on the spread of the COVID-19 pandemic: A review of findings, and statistical and modelling techniques. *Prog. Phys. Geogr.* (2020). doi:10.1177/0309133320946302

2.  WHO. *Coronavirus disease (COVID-19).* (2020).

3.  Vellingiri, B. *et al.* COVID-19: A promising cure for the global panic. *Science of the Total Environment* (2020). doi:10.1016/j.scitotenv.2020.138277

4.  Chan, K. H. *et al.* The effects of temperature and relative humidity on the viability of the SARS coronavirus. *Adv. Virol.* (2011). doi:10.1155/2011/734690

5.  Chin, A. W. H. *et al.* Stability of SARS-CoV-2 in different environmental conditions. *The Lancet Microbe* (2020). doi:10.1016/s2666-5247(20)30003-3

6.  Araújo, M. B., Mestre, F. & Naimi, B. Ecological and epidemiological models are both useful for SARS-CoV-2. *Nature Ecology and Evolution* (2020). doi:10.1038/s41559-020-1246-y

7.  Carlson, C. J., Chipperfield, J. D., Benito, B. M., Telford, R. J. & O'Hara, R. B. Species distribution models are inappropriate for COVID-19. *Nature Ecology and Evolution* (2020). doi:10.1038/s41559-020-1212-8

8.  Rohde, K. Latitudinal gradients in species diversity and Rapoport's rule revisited: A review of recent work and what can parasites teach us about the causes of the gradients? *Ecography (Cop.).* **22**, 593–613 (1999).

9.  Kindlmann, P., Schödelbauerová, I. & Dixon, A. F. G. Inverse latitudinal gradients in species diversity. *Scaling Biodivers.* 246–257 (2012). doi:10.1017/cbo9780511814938.014

10. Guernier, V., Hochberg, M. E. & Guégan, J. F. Ecology drives the worldwide distribution of human diseases. *PLoS Biol.* **2**, 740–746 (2004).

11. Lauring, A. S. & Andino, R. Quasispecies theory and the behavior of RNA viruses. *PLoS Pathogens* (2010). doi:10.1371/journal.ppat.1001005

12. Racevska, E. Natural Selection. in *Encyclopedia of Animal Cognition and Behavior* (eds. Vonk, J. & Shackelford, T.) 1–14 (Springer International Publishing, 2018). doi:10.1007/978-3-319-47829-6_542-1

13. Fehr, A. R. & Perlman, S. Coronaviruses: An overview of their replication and pathogenesis. in *Coronaviruses: Methods and Protocols* 1–23 (2015). doi:10.1007/978-1-4939-2438-7_1

14. Li, C. K. *et al.* T Cell Responses to Whole SARS Coronavirus in Humans. *J. Immunol.* **181**, 5490–5500 (2008).

15. Li, F., Li, W., Farzan, M. & Harrison, S. C. Structural biology: Structure of SARS coronavirus spike receptor-binding domain complexed with receptor. *Science (80-. ).* **309**, 1864–1868 (2005).

16. Liu, J. *et al.* The Membrane Protein of Severe Acute Respiratory Syndrome

577     Coronavirus Acts as a Dominant Immunogen Revealed by a Clustering Region
578     of Novel Functionally and Structurally Defined Cytotoxic T-Lymphocyte
579     Epitopes. *J. Infect. Dis.* **202**, 1171–1180 (2010).

580   17.   Ruch, T. R. & Machamer, C. E. The coronavirus E protein: Assembly and
581     beyond. *Viruses* (2012). doi:10.3390/v4030363

582   18.   Chang, C. K., Hou, M. H., Chang, C. F., Hsiao, C. D. & Huang, T. H. The
583     SARS coronavirus nucleocapsid protein - Forms and functions. *Antiviral*
584     *Research* (2014). doi:10.1016/j.antiviral.2013.12.009

585   19.   Allaby, M. *Atmosphere. A scientific history of air, weather, and climate.*
586     *Advances in Space Research* (2009).

587   20.   Peel, M. C., Finlayson, B. L. & McMahon, T. A. Updated world map of the
588     Köppen-Geiger climate classification. *Hydrol. Earth Syst. Sci.* **11**, 1633–1644
589     (2007).

590   21.   Wu, F. *et al.* A new coronavirus associated with human respiratory disease in
591     China. *Nature* (2020). doi:10.1038/s41586-020-2008-3

592   22.   Harcourt, B. H. *et al.* Identification of Severe Acute Respiratory Syndrome
593     Coronavirus Replicase Products and Characterization of Papain-Like Protease
594     Activity. *J. Virol.* (2004). doi:10.1128/jvi.78.24.13600-13612.2004

595   23.   Snijder, E. J., Decroly, E. & Ziebuhr, J. The Nonstructural Proteins Directing
596     Coronavirus RNA Synthesis and Processing. in *Advances in Virus Research*
597     (2016). doi:10.1016/bs.aivir.2016.08.008

598   24.   Zhang, L. *et al.* The D614G mutation in the SARS-CoV-2 spike protein reduces
599     S1 shedding and increases infectivity. *bioRxiv* 2020.06.12.148726 (2020).
600     doi:10.1101/2020.06.12.148726

601   25.   Korber, B. *et al.* Spike mutation pipeline reveals the emergence of a more
602     transmissible form of SARS-CoV-2. *bioRxiv* (2020).
603     doi:10.1101/2020.04.29.069054

604   26.   Junior, I. J. M. *et al.* The global population of SARS-CoV-2 is composed of six
605     major subtypes. *bioRxiv* (2020). doi:10.1101/2020.04.14.040782

606   27.   Becerra-Flores, M. & Cardozo, T. SARS-CoV-2 viral spike G614 mutation
607     exhibits higher case fatality rate. *Int. J. Clin. Pract.* (2020).
608     doi:10.1111/ijcp.13525

609   28.   WHO | SARS (Severe Acute Respiratory Syndrome). Available at:
610     https://www.who.int/ith/diseases/sars/en/. (Accessed: 10th September 2020)

611   29.   Boël, G. *et al.* Codon influence on protein expression in E. coli correlates with
612     mRNA levels. *Nature* **529**, 358–363 (2016).

613   30.   Tsai, C. J. *et al.* Synonymous Mutations and Ribosome Stalling Can Lead to
614     Altered Folding Pathways and Distinct Minima. *Journal of Molecular Biology*
615     **383**, 281–291 (2008).

616   31.   Shabalina, S. A., Ogurtsov, A. Y. & Spiridonov, N. A. A periodic pattern of
617     mRNA secondary structure created by the genetic code. *Nucleic Acids Res.*

618   (2006). doi:10.1093/nar/gkl287

619 32. Follis, K. E., York, J. & Nunberg, J. H. Furin cleavage of the SARS coronavirus
620   spike glycoprotein enhances cell-cell fusion but does not affect virion entry.
621   *Virology* **350**, 358–369 (2006).

622 33. Forster, P., Forster, L., Renfrew, C. & Forster, M. Phylogenetic network
623   analysis of SARS-CoV-2 genomes. *Proc. Natl. Acad. Sci. U. S. A.* (2020).
624   doi:10.1073/pnas.2004999117

625 34. Coronavirus Update (Live): 8,522,724 Cases and 453,714 Deaths from
626   COVID-19 Virus Pandemic - Worldometer.

627 35. Kumar, S., Stecher, G., Li, M., Knyaz, C. & Tamura, K. MEGA X: Molecular
628   evolutionary genetics analysis across computing platforms. *Mol. Biol. Evol.* **35**,
629   1547–1549 (2018).

630 36. Nei, M. & Saitou, N. The neighbor-joining method: a new method for reco...
631   [Mol Biol Evol. 1987] - PubMed result. *Mol Biol Evol* 406–425 (1987).

632 37. Tamura, K., Nei, M. & Kumar, S. Prospects for inferring very large phylogenies
633   by using the neighbor-joining method. *Proc. Natl. Acad. Sci. U. S. A.* **101**,
634   11030–11035 (2004).

635 38. EarthExplorer.

636 39. Travel Weather Averages (Weatherbase).

637 40. Climate data for cities worldwide - Climate-Data.org.

638 41. Herbei, M., Ciolac, V., Smuleac, A. & Ciolac, L. Georeferencing of
639   Topographical Maps Using the Software ArcGIS. *Res. J. Agric. Sci.* **42**, 595–
640   606 (2010).

641

## Authors Contribution

Phylogenetic study is carried out by P.B. GIS study & Koppen's climate map
interpretations is done by P.C.A. Worldometer data analysis is carried out by both
the authors. Both authors have written, reviewed & edited the manuscript.

## Acknowledgement

25

650    thank Prof. Raghavan Varadarajan, Prof. Raman Sukumar, Dr. Teena Jangid &

651    Chetankumar Jalihal of Indian Institute of Science for proofreading the article.

652    ## Conflict of Interest

653    Authors declare no conflict of interest.

654

655    # TABLE

656

657    **Table 1: SNPs representing virus cluster & their distribution across varied**
658    **climates.**

| Virus cluster | Nucleotide mutation | Amino acid mutation | Gene | Climate Zone | KCT | KC |
|---|---|---|---|---|---|---|
| 1-22 | C8782T | S2839 | ORF1a | NTZ | Cfa | C |
| | T28144C | L84S | ORF8 | | | |
| 5-6 | C29095T | F274 | N | NTZ | Cfa | C |
| 8-9 | T9477A | F3071Y | ORF1a | NTZ, TZ | Mix | C-A |
| | G25979T | G196V | ORF3a | | | |
| | C28657T | D128 | N | | | |
| | C28863T | S197L | N | | | |
| 10-17 | C18060T | L1431 | ORF1b | NTZ | Cfa-Cfb | C |
| 12-17 | A17858G | Y1364C | ORF1b | NTZ | Cfa-Cfb | C |
| 13-17 | C17747T | P1327L | ORF1b | NTZ | Cfa | C |
| 20-22 | C24034T | N824 | S | NTZ | Cfa | C |
| | T26729C | A69 | M | | | |
| | G28077C | V62L | ORF8 | | | |
| 21-22 | T490A | D75 | ORF1a | NTZ, NTSZ | Cfa | C |
| | C3177T | P971L | ORF1a | | | |
| | T18736C | F1657L | ORF1b | | | |
| 23-25 | C6312A | T2016K | ORF1a | NTZ, TZ, | Mix | Mix |
| | C13730T | L4489 | ORF1a | SSTZ | | |
| | C23929T | Y789 | S | | | |
| | C28311T | P13L | N | | | |
| 28-32 | G1397A | D392G | ORF1a | NTZ | Mix | Mix |
| | T28688C | L139 | N | | | |
| 33-43 | G26144T | G251V | ORF3a | NTZ | Cfa-Cfb | C |
| 37-39 | A2480G | I739V | ORF1a | NTZ | Mix | Mix |
| | C2558T | P765S | ORF1a | | | |
| 37-43 | C14805T | Y346 | ORF1b | NTZ | Cfa | C |
| 42-43 | T17247C | R1160 | ORF1b | NTZ | Cfb | B |
| 44-47 | G1440A | G392D | ORF1a | NTZ | Cfb | C |
| | G2891A | A876T | ORF1a | | | |
| 58-61 | C15324T | N519 | ORF1b | NTZ | Cfa-Dfb | C-D |
| 59-176 | C3037T | F924 | ORF1a | NTZ | Cfa-Cfb-Dfb-Aw | C-D -A |
| | A23403G | D614G | S | | | |
| | C14408T | P214L | ORF1b | | | |
| 59-125, | C241T | | Leader | NTZ | Cfa-Cfb-Dfa-Dfb | C-D-A |
| 127-176 | C241T | | seq. | | | C-D |

26

| 66-68 | A26530G | D3G | M | NTZ | Cfc-Dfb | C-D |
|---|---|---|---|---|---|---|
| 70-71 | G4201T | M1312I | ORF1a | NTZ | Cfa-Dwc | C-D |
| | C26527T | A2V | M | | | |
| 80-115 | G28881A | R203K | N | NTZ | | |
| | G28882A | R203K | N | | | |
| | G28883C | G204R | N | | Cfa-Cfb-Dfb-Aw | C-D-A |
| 86-87 | C27046T | T175M | M | NTZ | Cfa-Dfb | C-D |
| 88-89 | C3373A | D1036E | ORF1a | NTZ | Dfb-Cfb | C-D |
| 105-107 | T29148C | I292T | N | TZ, STZ | Cfa-Aw | C-A |
| 106-107 | A27299C | I33T | ORF6 | NTZ, TZ | Cfa-Aw | C-A |
| 108-111 | C313T | L16 | ORF1a | NTZ, TZ | Cfa-Cfb-Aw | C |
| 113-115 | C4002T | T1246I | ORF1a | STZ | Cfa-Cfb-Am | C-A |
| | G10097A | G3278S | ORF1a | | | |
| | C13536T | T4424I | ORF1a | | | |
| | C23731T | T723 | S | | | |
| 116-125 | A20268G | L2167 | ORF1b | NTZ | Cfa-Cfb-Dfa-Dfb | C-D |
| 126-176 | G25563T | Q57H | ORF3a | NTZ | Cfa-Cfb-Dfa-Dfb | C-D |
| 126-130 | C18877T | L1704 | ORF1b | NTZ | Cfa-Dfa-Dcb | C |
| 131-135 | C2416T | Y717 | ORF1a | NTZ | Cfa-Dfa-Aw | D |
| 136-176 | C1059T | T265I | ORF1a | NTZ | Cfa-Cfb-Dfa-Dfb | C-D |
| 138-139 | C18998T | A1744V | ORF1b | NTZ, TZ | Cfa-Am | C-A |
| | G29540A | | | | | |
| 138-141 | C11916T | S3884L | ORF1a | NTZ | Cfa-Csb-Am | C |
| 143-147 | C27964T | S24L | ORF8 | NTZ | Cfa-Cfb-Dfa-Dfb | C-D |
| 148-149 | C11224T | V3653 | ORF1a | NTZ | Dfa-Dfb | D |
| 157-159 | G29553A | | | NTZ | Cfa | C |

659

660  **NOTE:** Virus clusters are named by Strain ID as depicted on the tree (Supplementary table S1 & S2).
661  Genomic coordinates in this study is based on reference genome[21]. The SNP positions are based on
662  the reference genome. Nucleotide T represents nucleotide U in the SARS-CoV-2 RNA genome.
663  Mutation at the protein level is not mentioned for the SNPs arising in the non-coding region. The
664  amino acid position numbering is according to its position within the specified gene (CDS). In Climate
665  zone column we have mentioned the major climate zone for the corresponding virus cluster[19]. KCT is
666  Koppen's Climate Type & KC is Koppen's Climate columns display the main Koppen's climate in which
667  most of the virus isolates of the corresponding virus cluster lie. 'Mix' implies no particular climate type
668  is favored[20].

669

670

671

672

673

674

675

676

677

27

678

**(a)** **(b)**



NTZ (30°N - 66.5°N)    NSTZ (23.5°N - 30°N )

TZ (23.5°N - 23.5°S)    STZ (30°S - 66.6°S )

30°N - 37°N    37°N - 44°N    44°N - 51°N

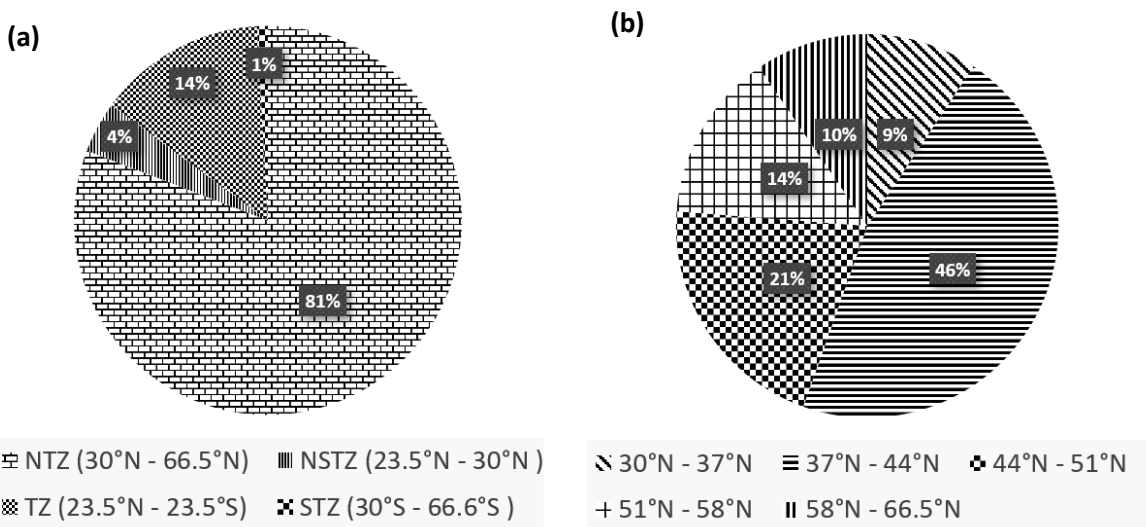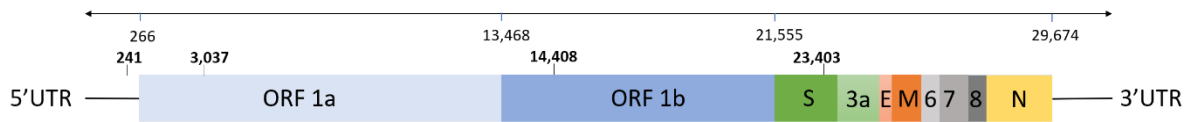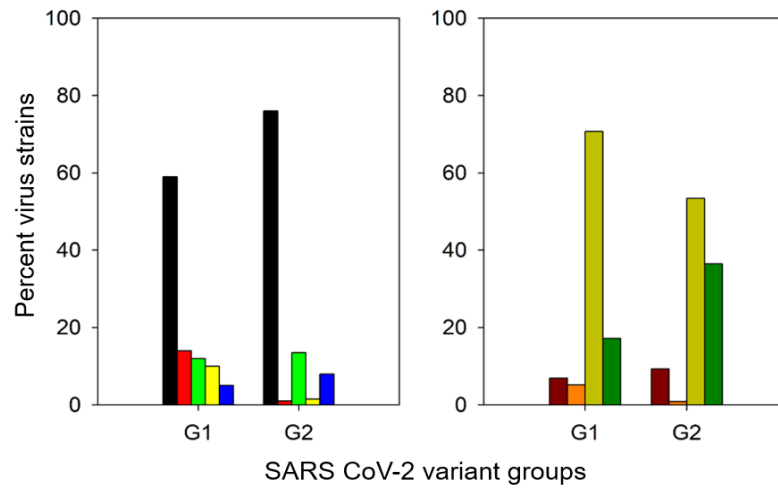51°N - 58°N    58°N - 66.5°N

679

680

681    **Figure1: Distribution of COVID-19 cases across different climate zones based on latitudes.** (a) Area
682    of the pie-chart covered by a climate zone is proportional to the percentage of COVID-19 cases
683    occurring in their respective climate zones as depicted by black squares. The percentage of COVID-19
684    cases for NFZ & SSTZ is extremely low, therefore, it is not mentioned in the pie-chart. (b) The North
685    Temperate Zone is divided into an interval of 7° latitude. The area of the pie-chart covered is directly
686    proportional to the percentage of COVID-19 cases occurring in their respective latitude range as
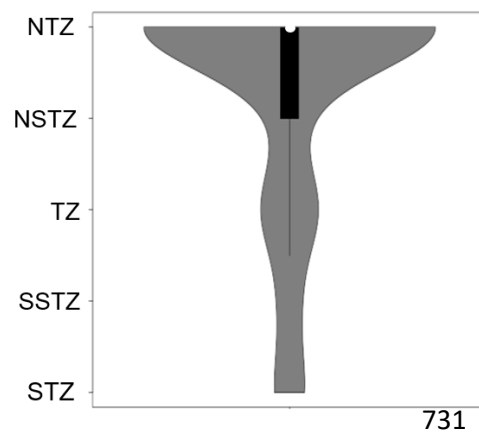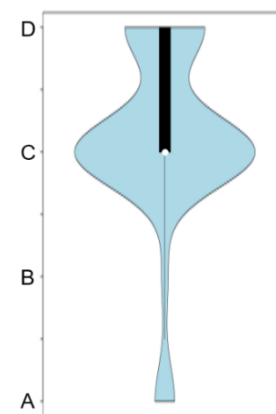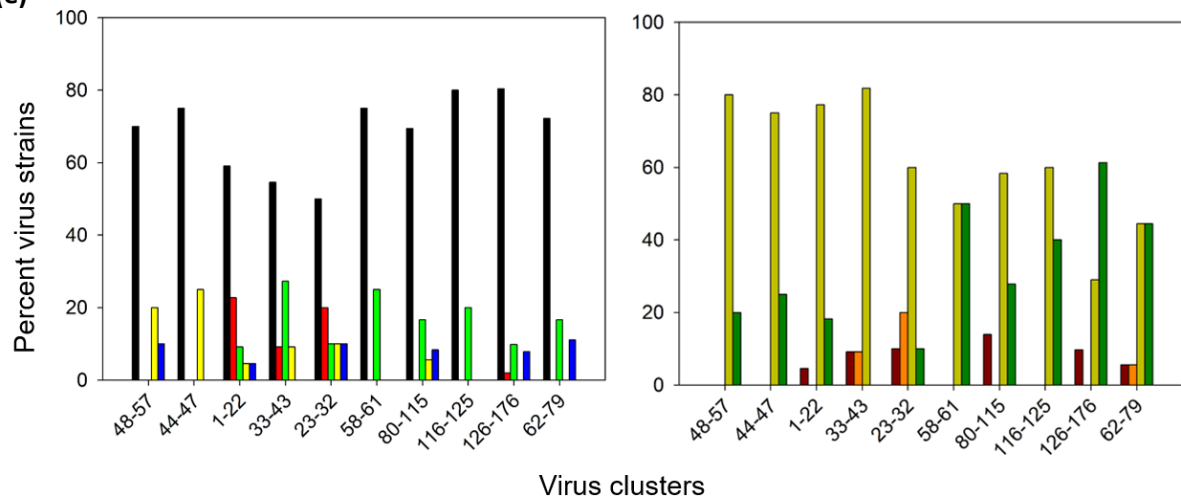687    depicted in black squares.

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706 **Figure2: Phylogenetic network divides 176 SARS-CoV-2 strains into two variant groups.** Largely, the
707 left side of the tree (1 to 58) constitute the G1 group & the right side of the tree constitutes the G2
708 group (59 to 176). Branch length is proportional to the genomic relatedness of the viral isolates.
709 Closely related virus isolates comprise the same SNP with respect to the reference genome (Strain ID:
710 50) & form a cluster. The evolutionary history of 176 taxa was inferred using the Neighbor-Joining
711 method[36] (500 bootstrap tests). A total of 29408 positions were analysed with nucleotide position
712 numbering according to the reference sequence[21].

29

**(a)**

713

**(b)**

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

**(c)**

**(d)**

731

**(e)**

732

733

734

735

736

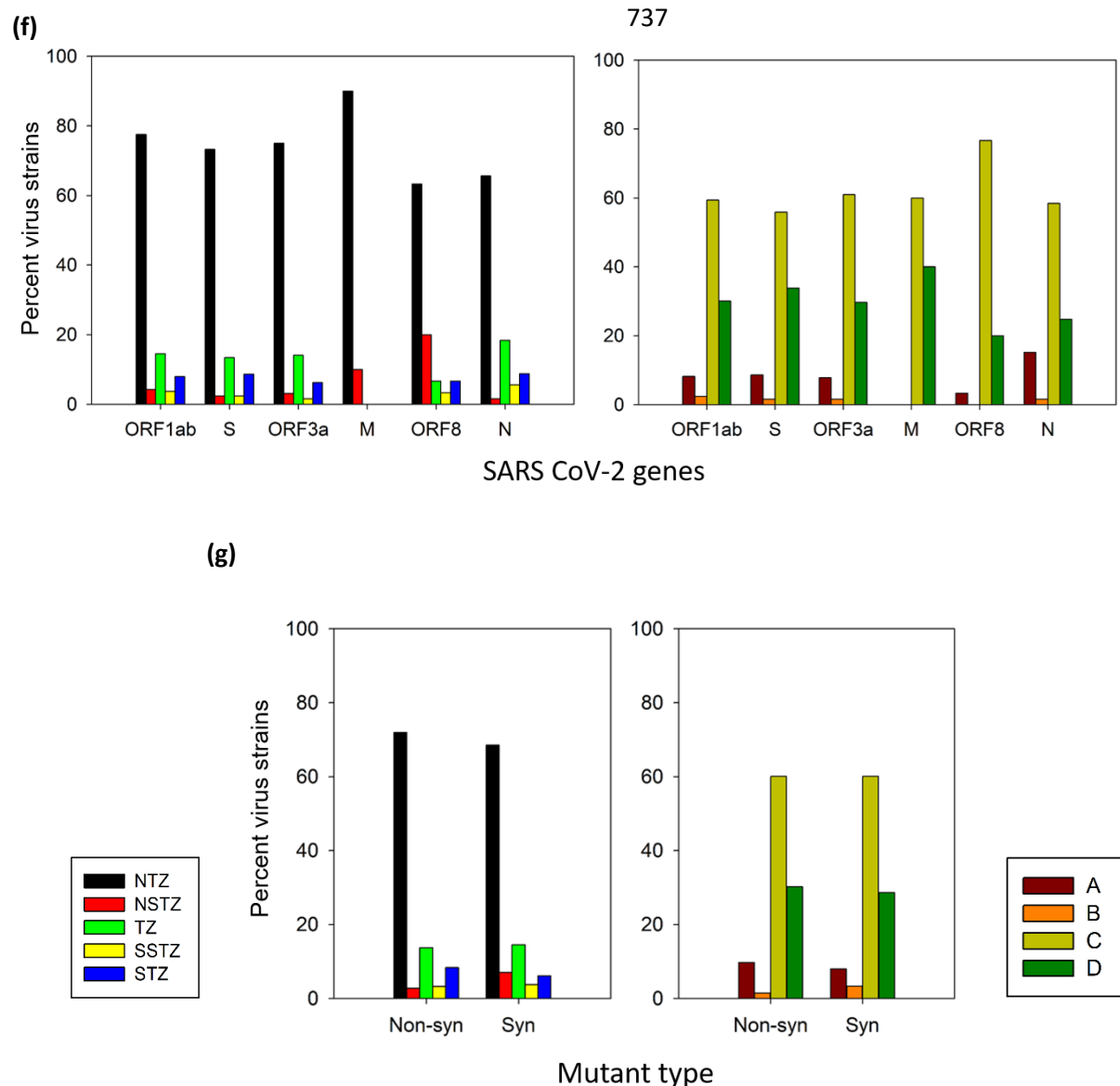737

738

739

740

741

742

743

744



745

746

**Figure3: Molecular phylogeny analysis to infer genomic similarities of SARS-CoV-2 & their distribution across different climate zones[19] & Koppen's climate types[20].** (a) Genomic architecture of SARS-CoV-2 genome highlighting four positions, substitutions on these positions enabled evolution of G1 into G2. (b, e-g) Strains found within a virus cluster (as shown in the phylogenetic tree & mentioned in Table 1) were analysed for significant mutations that may have arisen due to climatic pressure. Hence, percentage of such virus strains is plotted according to the geographical location of the climate zone from where they were isolated. The height of the bar is proportional to percent virus strain occurring in the specified condition i.e., labelled on the x-axis. Box in the left panel consist of color code for each climate zone & box in the right panel consist of color code for Koppen's climate. Left panel shows distribution of percent virus strains in different climate zones & right panel shows distribution of percent virus strain in Koppen's climate (b) Percent virus strains prevailing in different climate zones, stratified by SARS-CoV-2 variant groups. Width of curves of violin plot is proportional to the number of SARS-CoV-2 strains (n=176) in varied (c) climate zones & (d) Koppen's climate. (e) Abiotic factors influencing evolutionary dynamics of phylogenetic virus clusters. (f) Percent of virus strains with high frequency SNPs in each gene. (g) Type of mutation i.e. non-synonymous or synonymous exhibited by viruses.
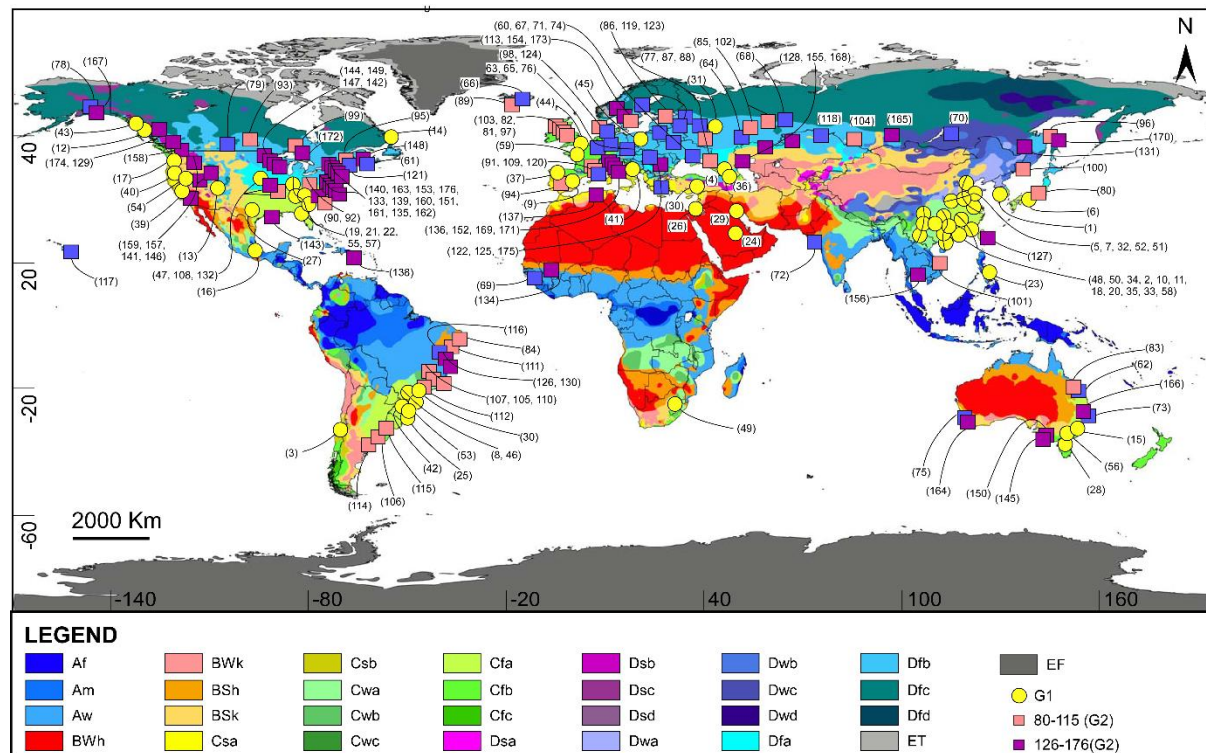
**Figure4: Global distribution of SARS-CoV-2 strains on the Gieger-Koppen's map displaying different climate types**[20]. Each strain is labelled as per the strain ID (1 to 176) within parenthesis. The G1 strains were symbolized as 'Yellow-circle', & G2 as 'Square', pink square denotes strain clusters (80-115) stable across C, D & A climate, purple square represents strain cluster (126-176) stable majorly in D climate, the remaining G2 strains (blue squares) are stable across C & D climate. Standard Koppen's climate-type symbols are mentioned in the legend, the criteria for distinguishing these climate types is mentioned in Table S3. Table S4 contains full form of these symbols. All symbols with initials 'A' (Af, Am, Aw) are of tropical climate, initials with 'B' belong to desert climate, 'C' to temperate & 'D' to cold & 'E' to polar climate. The shades of blue on the map, in North America & Russia belongs to D climate. Shades of yellow & green belongs to C climate, shades of red, orange & pink belongs to Desert climate.
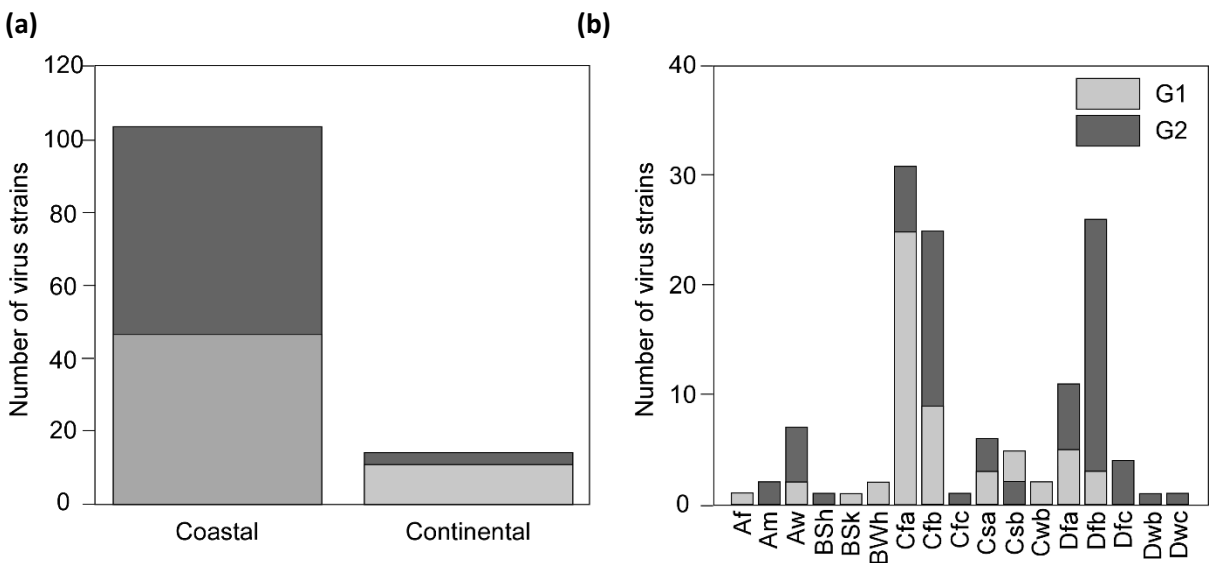
775



776

777 **Figure5: Global distribution of SARS-CoV-2 strains (n=176)** (a) in the coastal & continental region (b)
778 & in different Koppen's climate types[20]. Number of virus strains in G1 population is represented by
779 light grey color & of virus strains in G2 population is represented by dark grey color.

780

781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807

808

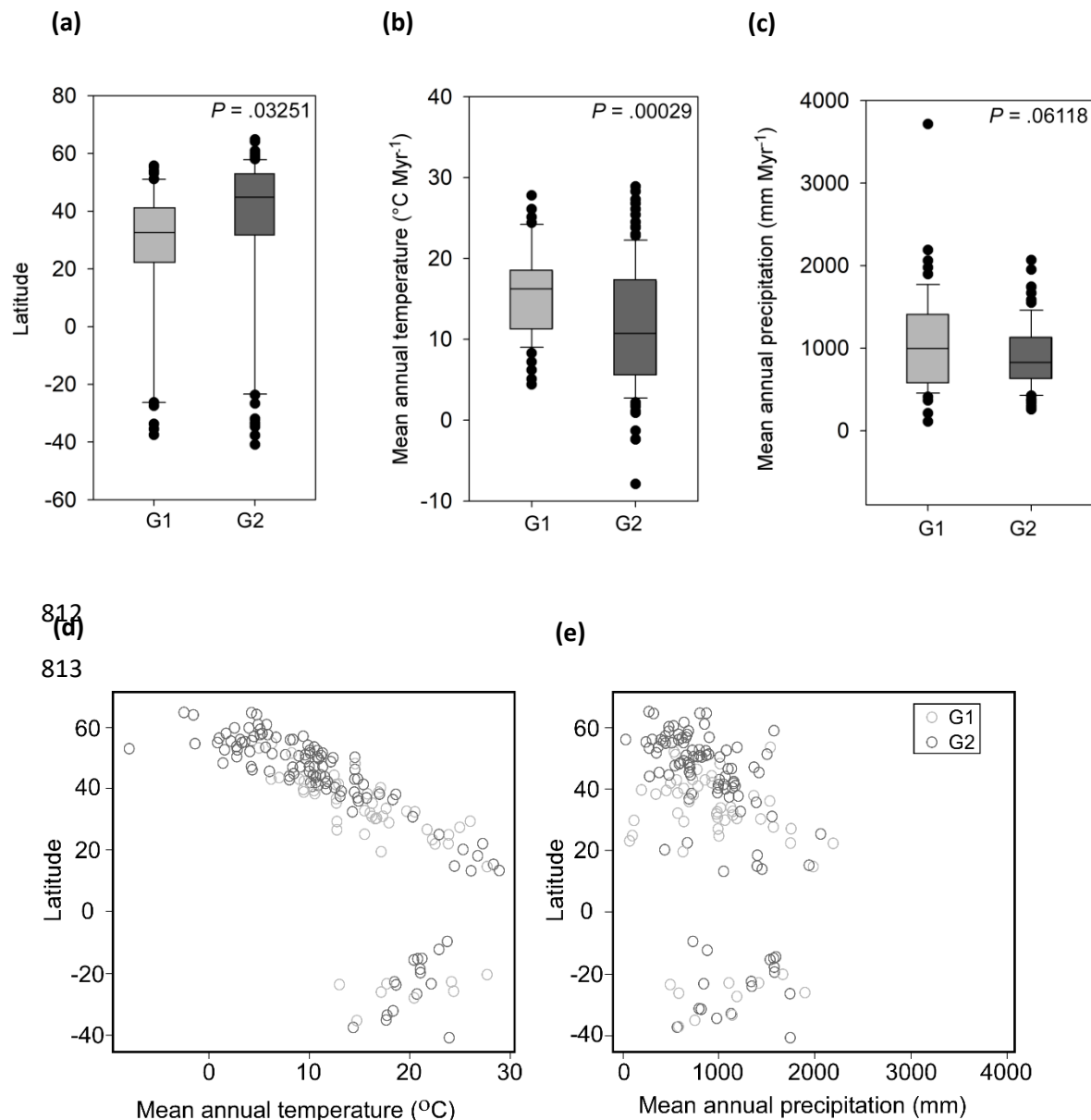809

810



811

812
813



822

**Figure6: Comparative analysis of different climatic parameters such as latitude, temperature & precipitation between G1 & G2 variant groups.** (a) Positive values represent the latitude range falling in Northern Hemisphere & negative values represent latitude range falling in Southern Hemisphere. The G2 strains preferentially occur towards the higher latitudes than G1 ($P$=.032; 95% CI 17.12-31.12 for G1; 95% CI 28.67-68.06 for G2). (b) The mean annual temperature of G2 is significantly lower than the G1 strains ($P$<.001; 95% CI 17.32-14.32 for G1; 95% CI 13.02-10.33 for G2) (c) Mean annual precipitation of G1 & G2 strains is nearly same ($P$=.061; 95% CI 1207.16-886.75 for G1; 95% CI 966.91-826.37 for G2). (a-c) Black horizontal line in the middle of the box is median, upper & lower limits of the box indicate first & third quartile. Black dots represent outliers. P values is based on one-way ANOVA. Scatter plot for (d) latitude & annual temperature & (e) latitude & precipitation for each SARS-CoV-2 strain (n=176) belonging to G1 group (n=58, shown in light grey) & G2 group (n=118, shown in dark grey).

835