

1 **Landscape of SARS-CoV-2 genomic surveillance, public availability extent of genomic**
2 **data, and epidemic shaped by variants: a global descriptive study**

3

4 Zhiyuan Chen¹, Andrew S. Azman^{2,3}, Xinhua Chen¹, Junyi Zou¹, Yuyang Tian¹, Ruijia Sun¹,
5 Xiangyanyu Xu¹, Yani Wu¹, Wanying Lu¹, Shijia Ge⁴, Zeyao Zhao¹, Juan Yang¹, Daniel T. Leung^{5,6},
6 Daryl B. Domman⁷, and Hongjie Yu^{1,4,8}

7

8 **Affiliations**

- 9 1. School of Public Health, Fudan University, Key Laboratory of Public Health Safety, Ministry of
10 Education, Shanghai, China
- 11 2. Department of Epidemiology, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD,
12 USA
- 13 3. Institute of Global Health, Faculty of Medicine, University of Geneva, Switzerland
- 14 4. Department of Infectious Diseases, Huashan Hospital, Fudan University, Shanghai, China
- 15 5. Division of Infectious Diseases, University of Utah School of Medicine, Salt Lake City, UT, USA
- 16 6. Division of Microbiology & Immunology, University of Utah School of Medicine, Salt Lake City,
17 UT, USA
- 18 7. Center for Global Health, Department of Internal Medicine, University of New Mexico Health
19 Sciences Center, New Mexico, USA
- 20 8. Shanghai Institute of Infectious Disease and Biosecurity, Fudan University, Shanghai, China

21

22 Corresponding authors: Hongjie Yu, School of Public Health, Fudan University, Shanghai 200032,
23 China; E-mail: yhj@fudan.edu.cn

24

25 Word count (abstract): 286

26 Word count (main text): 3,534

27

28 **Disclaimer:** The views expressed are those of the authors and do not necessarily represent the
29 institutions with which the authors are affiliated.

30 **Abstract**

31 **Background**

32 Genomic surveillance has shaped our understanding of SARS-CoV-2 variants, which have
33 proliferated globally in 2021. Characterizing global genomic surveillance, sequencing coverage, the
34 extent of publicly available genomic data coupled with traditional epidemiologic data can provide
35 evidence to inform SARS-CoV-2 surveillance and control strategies.

36

37 **Methods**

38 We collected country-specific data on SARS-CoV-2 genomic surveillance, sequencing capabilities,
39 public genomic data, and aggregated publicly available variant data. We divided countries into three
40 levels of genomic surveillance and sequencing availability based on predefined criteria. We
41 downloaded the merged and deduplicated SARS-CoV-2 sequences from multiple public repositories,
42 and used different proxies to estimate the sequencing coverage and public availability extent of
43 genomic data, in addition to describing the global dissemination of variants.

44

45 **Findings**

46 Since the start of 2021, the COVID-19 global epidemic clearly featured increasing circulation of
47 Alpha, which was rapidly replaced by the Delta variant starting around May 2021 and reaching a
48 global prevalence of 96.6% at the end of July 2021. SARS-CoV-2 genomic surveillance and
49 sequencing availability varied markedly across countries, with 63 countries performing routine
50 genomic surveillance and 79 countries with high availability of SARS-CoV-2 sequencing. Less than
51 3.5% of confirmed SARS-CoV-2 infections were sequenced globally since September 2020, with the
52 lowest sequencing coverage in the WHO regions of Eastern Mediterranean, South East Asia, and
53 Africa. Across different variants, 28-52% of countries with explicit reporting on variants shared less
54 than half of their variant sequences in public repositories. More than 60% of demographic and 95% of
55 clinical data were absent in GISAID metadata accompanying sequences.

56

57 **Interpretation**

58 Our findings indicated an urgent need to expand sequencing capacity of virus isolates, enhance the
59 sharing of sequences, the standardization of metadata files, and supportive networks for countries
60 with no sequencing capability.

61

62 **Keywords**

63 Genomic Surveillance, Sequencing, SARS-CoV-2 Variants, Data Sharing

64

65 **Funding**

66 Key Program of the National Natural Science Foundation of China, the US National Institutes of
67 Health.

68 **Research in context**

69 **Evidence before this study**

70 On September 3, 2021, we searched PubMed for articles in any language published after January 1,
71 2020, using the following search terms: (“COVID-19” OR “SARS-CoV-2”) AND (“Global” OR
72 “Region”) AND (“genomic surveillance” OR “sequencing” OR “spread”). Among 43 papers
73 identified, few papers discussed the global diversity in genomic surveillance, sequencing, public
74 availability of genomic data, as well as the global spread of SARS-CoV-2 variants. A paper from
75 Furuse employed the publicly GISAID data to evaluate the SARS-CoV-2 sequencing effort by
76 country from the perspectives of "fraction", "timeliness", and "openness". Another viewpoint paper by
77 Case Western Reserve University’s team discussed the impediments of genomic surveillance in
78 several countries during the COVID-19 pandemic. The paper as reported by Campbell and colleagues
79 used the GISAID data to present the global spread and estimated transmissibility of recently emerged
80 SARS-CoV-2 variants. We also found several studies that reported the country-level genomic
81 surveillance and spread of variants. To our knowledge, no research has quantitatively depicted the
82 global SARS-CoV-2 genomic surveillance, sequencing ability, and public availability extent of
83 genomic data.

84

85 **Added value of this study**

86 This study collected country-specific data on SARS-CoV-2 genomic surveillance, sequencing
87 capabilities, public genomic data, and aggregated publicly available variant data as of 20 August 2021.
88 We found that genomic surveillance strategies and sequencing availability is globally diverse. Less
89 than 3.5% of confirmed SARS-CoV-2 infections were sequenced globally since September 2020. Our
90 analysis of publicly deposited SARS-CoV-2 sequences and officially reported number of variants
91 implied that the public availability extent of genomic data is low in some countries, and more than 60%
92 of demographic and 95% of clinical data were absent in GISAID metadata accompanying sequences.
93 We also described the pandemic dynamics shaped by VOCs.

94

95 **Implications of all the available evidence**

96 Our study provides a landscape for global sequencing coverage and public availability extent of
97 sequences, as well as the evidence for rapid spread of SRAS-CoV-2 variants. The pervasive spread of
98 Alpha and Delta variants further highlights the threat of SARS-CoV-2 mutations despite the
99 availability of vaccines in many countries. It raised an urgent need to do more work on defining the
100 ideal sampling schemes for different purposes (e.g., identifying new variants) with an additional call
101 to share these data in public repositories to allow for further rapid scientific discovery.

102 **Introduction**

103 Following the first pandemic wave of coronavirus disease 2019 (COVID-19), the emergence and
104 dissemination of SARS-CoV-2 variants have resulted in new waves of infections across the globe in
105 2021. Some SARS-CoV-2 variants disappeared immediately, while others characterized by several
106 key mutations adapted well, enabling their rapid spread¹. WHO has designated four Variants of
107 Concern (VOCs) associated with increased transmissibility and various extents of immune escape²⁻⁴,
108 namely the Alpha, Beta, Gamma, and Delta variants, first detected in the UK, South Africa, Brazil,
109 and India, respectively⁵. Specifically, the Delta variant is highly transmissible, with an estimated
110 transmissibility increase of 40-60% compared with Alpha variant⁶⁻⁹, while the Beta variant has been
111 shown to have the highest reduction in neutralization activity whether from natural infection or
112 vaccination¹⁰, both reflected in lower vaccine efficacy or effectiveness¹¹⁻¹³.

113
114 The identification and classification of SARS-CoV-2 variants mainly relied on partial or whole
115 genome sequencing, although PCR assays have been used to identify specific features relatively
116 unique in specific variants, like spike gene target failure (SGTF)¹⁴. Since the first SARS-CoV-2
117 sequence was published in January 2020¹⁵, the unprecedented rate of genome data generation was far
118 greater than any other pathogen¹⁶, with 3.1 million genomes deposited in Global Initiative on Sharing
119 All Influenza Data (GISAID) through August 2021¹⁷. Genomic data has been vital to the early
120 detection of mutations and monitoring of virus evolution, as well as evaluating the degree of
121 similarities between circulating variants with vaccine strains, especially since the availability of
122 SARS-CoV-2 vaccines¹⁸.

123
124 Several studies have employed genomic data to examine the evolution and associated spread of
125 dominant variants in one country or one region, raising a claim of the rapidity of local transmission of
126 SARS-CoV-2 variants and urgency of genomic surveillance^{6,19-21}. However, the lack of representation
127 of genomic data from low and middle-income countries (LMICs) was most concerning in these
128 studies^{6,19-21}. Indeed, the impact of genome data is dependent on their quality, and the reliability and
129 accuracy of such data may influence the global community's ability to track the spread of variants in a
130 timely manner.

131
132 In this study, we aimed to investigate the global diversity of SARS-CoV-2 genomic surveillance and
133 the global sequencing coverage of confirmed cases and public availability extent of genomes. In
134 addition, we sought to map the global identification and spread of SARS-CoV-2 variants. This data
135 can provide evidence to better inform policy on SARS-CoV-2 surveillance.

136 **Methods**

137 **Data sources and collection**

138 Through extracting country-specific data from multiple publicly available sources, we built three
139 datasets of genomic surveillance, deposited genomic data in publicly repositories, and official
140 aggregated genomic data as of 20 August 2021.

141

142 *Dataset of genomic surveillance*

143 Each country's genomic surveillance strategy and sequencing availability was gathered from searches
144 of the websites of regional WHO, the country's Ministry of Health, CDC, local academic partners,
145 and official news, supplemented by a literature search (appendix). Data extracted included the overall
146 surveillance strategy, sequencing availability, target population, sampling method, diagnostic criteria,
147 and sequenced volume. Given that the surveillance strategy and density may change with time, we
148 only gathered information on the most recent surveillance strategy (as of 20 August 2021).

149

150 *Dataset of SARS-CoV-2 sequences in publicly repositories*

151 SARS-CoV-2 sequences along with related metadata file were downloaded from an online
152 coronavirus analysis platform from the National Genomics Data Center (NGDC)²², where has merged
153 and deduplicated sequences that deposited in GISAID²³, GenBank²⁴, National Genomics Data
154 Center²⁵, National Microbiology Data Center²⁶, and China National GeneBank²⁷. Detailed process of
155 integration and deduplication was previously reported²⁸. An acknowledgement table for those
156 contributing to this work is available in the Appendix.

157

158 *Official aggregated dataset*

159 To gain additional insights regarding the public availability (sharing) extent of genomic data of
160 SARS-CoV-2 variants, we extracted country-specific, variant-specific, and time-specific aggregated
161 data on the number of SARS-CoV-2 variant cases from official websites, using the same sources as
162 above, except for the literature source. The search was done by either directly locating to the official
163 website for each country or indirectly searching in search engines (Google, Bing, Baidu) by using the
164 terms "variant" and country name. To supplement the aggregated data of variants that we collected,
165 we also downloaded the aggregated data with a valid denominator (namely, the number of isolates
166 sequenced is reasonable) from the European Surveillance System (TESSy). The aggregated dataset
167 included country name, date of report or collection, new or cumulative numbers of different SARS-
168 CoV-2 variant cases and total sequenced cases, as well as the new number of confirmed cases.
169 COVID-19 epidemic data were derived from WHO²⁹. Considering that the diagnostic criteria of
170 SARS-CoV-2 variants varies in different countries, the general principle for collecting aggregated

171 data was to give priority to the results based on whole and partial genome sequencing instead of those
172 based on a PCR assay.

173

174 For countries noted by WHO as having had VOCs identified, but no data in publicly repositories, we
175 collected the information about when VOC/VOCs were first detected from country's Ministry of
176 Health and official media news. The search of media news was also done in search engine queries
177 (Google, Bing, Baidu) using combined terms "first" and "variant" and country name.

178

179 All data were entered into a structured database by a trained team (co-authors). All recorded data were
180 cross-checked by coauthors. Details of data sources and data completeness are shown in Appendix
181 (Table S1-S2, S4-S5).

182

183 **Data analysis**

184 We used the variant naming system proposed by WHO, where four VOCs and five Variants of
185 Interest (VOIs, included Eta, Iota, Kappa, Lambda, and Mu) had been designated as of 2 September
186 2021⁵. Our study focused on analyses on four VOCs in the 194 Member States of WHO. We did not
187 integrate data from the overseas territories into that country's data. Given most countries weekly
188 released aggregated data of variant, most of our analyses were performed on a weekly basis.

189

190 To characterize the global diversity of SARS-CoV-2 genomic surveillance, we classified the
191 surveillance strategy of each country into three categories: 1) routine genomic surveillance; 2) limited
192 routine genomic surveillance; and 3) no routine genomic surveillance. Routine genomic surveillance
193 was defined as conducting nationwide genomic sequencing, coupled with at least 150 specimens per
194 week or 10% of all samples sequenced³⁰. The classification of most African countries was from a pre-
195 defined version from African CDC³¹, while other countries were subsequently defined according to
196 the definition criteria (Table S3). As we could not identify information on the surveillance strategy for
197 some countries through public sources, we also classified three extra categories according to the
198 public availability/ability of genomic sequencing: 1) high availability; 2) moderate availability; and 3)
199 low availability (Table S3).

200

201 The process of data cleaning in publicly repositories and aggregated dataset is shown in the Appendix.
202 For data that was reported in aggregate, when the date of sample collection was not available, we
203 assumed a fixed three-week lag³² from sample collection to reporting unless other country-specific
204 information was available to inform this extrapolation.

205

206 The sequencing coverage was inferred by using the percent of cumulative positives sequenced as a
207 proxy, defined by the ratio of the number of isolates sequenced to the number of confirmed cases in
208 the same unit of time. Given that not all sequences will be uploaded to genomic repositories, we
209 analysed the public availability extent of genomic data by comparing the cumulative number of
210 variants between public repositories and official aggregated datasets. Since the Alpha variant had a
211 characteristic deletion of amino acids 69-70 in the SGTF that can be detected via a widely used PCR
212 assay³³, we performed this analysis across VOCs.

213

214 We plotted the earliest time when the first VOC specimen was identified in each country. The earliest
215 identification was defined by the earliest sampling time of sequences deposited in publicly
216 repositories. If a VOC was identified by WHO but no corresponding sequence in publicly repositories
217 for one country, we used the date collected from other sources. The sequences with sampling date
218 earlier than the earliest sample identified in the United Kingdom (for Alpha), South Africa (for Beta),
219 Brazil (for Gamma) and India (for Delta), respectively, were not used in analyses.

220

221 We also described the global and regional prevalence trend of variants. The prevalence of variant was
222 defined as the proportion of the variant number to the sequencing number in the same period. When
223 multiple data sources were available for one country, the most abundant dataset was chosen. For
224 example, if there was publicly available genomic data as well as an official aggregated dataset, the
225 priority was given to the one with highest reported sequenced numbers in a specific week.

226

227 Statistical significance was tested using the Chi-square test for comparing ratios of two groups, and
228 differences were considered statistically significant at P-value < 0.05. All statistical analyses and
229 visualizations were done using R (version 4.0.2).

230

231 **Results**

232 We classified genomic surveillance strategies for 105 countries, including 72.3% (34/47) of WHO-
233 defined Africa Region countries, 58.5% (31/53) of European Region countries, 61.9% (13/21) of
234 countries in the Eastern Mediterranean Region, 31.4% (11/35) of countries in Americas Region, 37.0%
235 (10/27) of countries in the Western Pacific Region, and 54.5% (6/11) of countries in the South East
236 Asia Region (Table S4). We downloaded a total of 3.1 million deduplicated SARS-CoV-2 sequence
237 samples from publicly repositories corresponding to samples collected between 1 December 2019 to
238 20 August 2021 in 163 countries. Additionally, we collected official aggregated data of variants from
239 55 countries, and extra data for the first identification of COVID-19 variants from 33 countries.

240

241 **The global diversity of SARS-CoV-2 genomic surveillance strategies and sequencing availability**

242 We observed marked geographical heterogeneity in genomic surveillance of SARS-CoV-2 across
243 countries. Globally, a total of 60.0% (63) countries had performed routine genomic surveillance, 26.7%
244 (28) countries implemented limited routine genomic surveillance, 13.3% (14) countries had no routine
245 genomic surveillance with the remaining countries (89) having no data on genomic surveillance
246 strategy identified (Figure 1A). Surveillance diversity across various countries was also reflected in
247 the context of target populations, sampling method, sequenced proportion, and diagnostic criteria
248 (Table S4). Specifically, 32 countries randomly selected or used all confirmed cases with sufficient
249 quality for sequencing, and 18 countries adopted PCR assay to screen or confirm variants. From the
250 regional perspective, limited or no routine genomic surveillance was common in the Eastern
251 Mediterranean (84.6%, 11/13) and Africa (61.8%, 21/34), followed by the South East Asia (50.0%,
252 3/6), Americas (36.4%, 4/11), Western Pacific (20.0%, 2/10), and Europe (3.2%, 1/31) (Figure 1A).
253 However, among the 167 Member States with data accessible, the availability of SARS-CoV-2
254 sequences was high in 79 countries, moderate in 80 countries, and low in 8 countries (Figure 1B).

255

256 **Sequencing coverage of SARS-CoV-2 confirmed cases**

257 The global volume of genomic data gradually increased over time until April 2021 (Figure 1C). The
258 European (59.4%) and Americas (33.2%) Regions uploaded the most sequences to public repositories,
259 with marked intra-region heterogeneity across countries, with a range from 3 (Saint Kitts and Nevis)
260 to 847,000 (United States) as of 20 August 2021 (Figure 1D).

261

262 Since September 2020, no more than 3.5% of global confirmed SARS-CoV-2 infections were
263 sequenced. In any week, the Africa, South East Asia and Eastern Mediterranean Regions (Figure 1E)
264 sequenced no more than 1.2% of cases. At the end of May 2021 (start of Delta widely spreading),
265 Europe had the highest sequenced proportion of 9.3%, followed by Western Pacific (2.1%), Americas
266 (1.4%), Africa (1.0%), South-East Asia (0.1%), and Eastern Mediterranean (0.07%) (Figure 1E). At
267 the country-level, higher rates of sequencing were observed in Iceland, New Zealand, Australia,
268 Denmark, Luxembourg, Finland, United Kingdom, and Norway, all of which had at least 10%
269 reported infections sequenced as of 20 August 2021. In addition, almost all countries in the Africa and
270 Eastern Mediterranean had sequenced less than 2.5% of confirmed infections, except for Gambia
271 (6.6%) and Mauritius (3.5%) (Figure 1F).

272

273 **Public availability extent of SARS-CoV-2 variant sequences**

274 The public availability extent of SARS-CoV-2 genomic data varied across variants and countries.
275 Overall, among countries with aggregated data on the number of variant infections, less than half of
276 sequences of Alpha, Beta, Gamma, and Delta were publicly available in 45.3% (24/53), 27.9%
277 (12/43), 33.3% (11/33), and 52.3% (23/44) countries, respectively (Figure 2). However, the result for

278 Alpha might be influenced by SGTF detected via PCR. The public availability extent of Delta variants
279 across countries ranged from 0.0% (Cyprus, Hungary, Iceland, Laos) to 92.9% (Denmark). At the
280 country level, low sharing proportions (less than 50%) across all VOCs was found in several countries,
281 including Austria, Cyprus, Greece, Hungary, Iceland, Philippines, Thailand and Senegal. For example,
282 the sharing proportion of Alpha, Beta, and Delta in Thailand was 5.2%, 13.6%, and 2.0%,
283 respectively, which indicated more than 85.0% genomic data of variants were not uploaded.

284

285 Moreover, incomplete metadata attached to GISAID sequences was common globally, with about two
286 thirds of sequences missing demographic information (age and sex), and more than 95% of that
287 missing clinical information (e.g., symptom history, clinical outcome, and vaccination status) (Table
288 S5). We found that high-income regions tend to have lower information completeness ($P < 0.0001$,
289 Chi-square test), especially in the European Region, where less than 25.0% and 3.0% of sequences
290 had demographic and clinical information, respectively. Besides, 92.8% sequences were reported by
291 subnational geographies, with a low proportion in Western Pacific (52.0%) and Eastern
292 Mediterranean (76.9%).

293

294 **Earliest identification of SARS-CoV-2 variants across regions**

295 Alpha was first identified in the Europe, then in Americas and South-East Asia in September-October
296 2020, followed by spread to Eastern Mediterranean, Africa, and Western Pacific in November 2020
297 (Figure 3A). The earliest publicly available sequenced Beta infection was sampled in Africa in May
298 2020, and subsequently identified in other regions (Figure 3B). Gamma variant detection remained
299 spatially constrained after it was first identified in Brazil (Figure 3C). After the first identification of
300 Delta in October 2020 in South-East Asia, global spread occurred after January 2021 (Figure 3D).

301

302 **Global and regional spread of SARS-CoV-2 variants**

303 The number of new VOC cases dramatically increased until April 2021, with a peak weekly value of
304 about 100,000 VOC cases sequenced in which most of them were Alpha variants (Figure 5A).

305 Subsequently, another peak of weekly new VOC case occurred in July 2021, but with a large amount
306 of Delta variants. The number of VOC cases may be an underestimate for the most recent weeks due
307 to a collection-to-report time delay. Notably, this increase was also accompanied by the increase in
308 the volume of new sequenced cases and new COVID-19 confirmed cases.

309

310 The global prevalence of reference (non-variant) strains fell into a low level of 0.6% in the period of
311 Jul-Aug 2021, compared with 13.4% of that in 2020 (Figure 4). Globally, the COVID-19 pandemic
312 was driven by the circulation of Alpha at the start of 2021, with an average prevalence of 50.6% in the
313 first quarter of 2021. Alpha variants continued to outcompete other strains in the second quarter of

314 2021, accounting for 60.0% of the contemporary lineages (Figure 4). However, the rapid global rise
315 of the Delta variant began in May 2021, reaching a global prevalence of nearly 96.6% at the end of
316 July 2021 (Figure 5B). In contrast, Beta and Gamma variants remained at low prevalence (Figure 4).
317 Additionally, the shifting of predominant variants from Alpha to Delta first occurred in South-East
318 Asia where the proportion of Delta exceeded 60.0% in April 2021 (Figure 5J).

319 Discussion

320 Our study characterized the global diversity of genomic surveillance strategies and sequencing
321 availability, sequenced coverage of SARS-CoV-2 cases, public availability extent of variant
322 sequences, as well as current epidemic situation of SARS-CoV-2 variants. We found that genomic
323 surveillance strategies were globally heterogenous, with limited or no routine surveillance among
324 many countries in the Africa and Eastern Mediterranean Regions. Our analysis of publicly deposited
325 SARS-CoV-2 sequences implied that the sequenced coverage is low in most countries, with a low
326 proportion of VOCs sequences shared to public repositories. The pervasive spread of Alpha and Delta
327 variants further highlights the threat of SARS-CoV-2 mutations despite the availability of vaccines in
328 many countries. Overall, we describe the pandemic dynamics shaped by VOCs and call for more
329 work on defining the ideal sampling schemes for different purposes (e.g., identifying new variants)
330 with an additional call to share these data in public repositories to allow for further rapid scientific
331 discovery.

332
333 The diversity of SARS-CoV-2 genomic surveillance between countries is associated with country-
334 specific priorities (e.g., surveillance objectives, targeted monitoring, or event-/risk-based sequencing)
335 and available resources. ECDC recommends population-based and/or targeted sampling strategies
336 (e.g., imported cases, cluster cases, and potential vaccine escapers) for genomic surveillance, which
337 could provide a more representative estimate of the relative prevalence of variants. Notably, several
338 countries, many of which are classified as low- or lower middle-income countries by the World Bank,
339 lack genomic surveillance data, likely due to limitations in infrastructure capacity and resources.
340 However, even some countries classified as high-income, have suffered from a slow and inconsistent
341 process of adopting genomics-based surveillance³⁴. Establishment of reference laboratories and
342 networks to provide sequencing services for countries without established sequencing capacity may
343 enable improved detection and tracking of emerging variants worldwide.

344
345 The detection of most variants relies on the full-length or partial genomic sequencing, but the
346 sequences only become available for the global community when the laboratories have established
347 sequencing capacity, willing to share, and legally allowed to upload them. The discrepancies in
348 sharing was observed in each region, which confirmed that some countries are sequencing but are not
349 uploading. However, our study observed a sharing extent of exceed 100% exists in some countries,
350 likely due to delays in the official reporting of sequencing results, or the incomplete official reporting
351 system. The timely sharing of those enables to adequately contextualize local data when looking at
352 introductions and examine transmission routes, as well as to look for sites of repeated mutations that
353 can guide laboratory work on characterizing those mutations effects on therapeutics and vaccine

354 efficacy. The underlying reasons why some countries didn't share might be related with the distrust
355 for publicly repositories in the concept of data security.

356

357 We found relatively low completeness of demographic and clinical characteristics in metadata
358 accompanying uploaded sequences. Our analysis revealed that high-income countries frequently did
359 not share demographic information. A possible reason for this is that these regions may having more
360 restrictive data privacy/laws preventing/discouraging the release of this information. Genomic data
361 coupled with those additional data can maximize the utility of genomic sequences in rapid scientific
362 discovery during this pandemic, which are valuable for in-depth epidemiological analyses to
363 characterize risk factors, clinical severity, and other public health risk of variants³⁵⁻³⁷. Therefore, it's
364 vital to optimize the sharing of information in a secure and trusted channel in the context of protect
365 patient anonymity and in accordance of local regulations³⁸. In addition, decreasing the lag between
366 sample collection to deposition of these sequences³⁹, including the timely sharing and standardizing
367 of metadata^{18,30,35}, may facilitate the design and development of treatment and prevention strategies
368 by policy-makers⁴⁰.

369

370 An important role of genomic surveillance is to investigate the spread and dynamics of SARS-CoV-2
371 variants. Amidst the emergence of different variants, the current dominance of the Delta variant
372 suggests that it may possess higher fitness than other variants, which might be associated with a
373 combination of higher virus load, and shorter incubation period and serial intervals⁴¹⁻⁴³. The decrease
374 in real-world vaccine effectiveness against Beta variants¹¹ and increasing breakthrough cases with
375 Delta variants⁴⁴ underlines the importance of determining the local or regional patterns of variant
376 spread, including the need to develop new or modified vaccines to achieve adequate protection⁴⁵.

377

378 Our results should be interpreted in view of several limitations. First, the lack of data from some
379 countries limited our global mapping. The data completeness and quality could be impacted by key
380 steps in the surveillance or reporting, including differences in diagnostic criteria, under-reporting,
381 delayed reporting, and reporting methods. The inconsistent diagnostic criteria of variants might cause
382 sampling bias, especially when adopting PCR assay to detect Alpha variant owing to its non-
383 specificity⁴⁶. We did an extensive search to collect multi-source data and chose the aggregated data
384 with a priority to sequencing results rather than PCR-screening results. Second, the analysis of global
385 and national spread could be biased as data from public repositories or aggregated dataset are not
386 always representative of the variants circulating in the regions, especially for the regions with
387 relatively limited sequencing capacity or with investigating outbreak-based events. Therefore, the
388 global prevalence of variants may be biased due to the uneven sequencing across the regions. Indeed,
389 it's difficult to obtain a truly representative and random sample, and how to understand these biases
390 will become important³⁶. Lastly, the detailed demographical, epidemiological and clinical information

391 about variant cases cannot be accessed, which limited our further epidemiological analysis about
392 variant spread.

393

394 In conclusion, our study provides a landscape for genomic surveillance, the global coverage of
395 sequencing and public availability of sequences, as well as the evidence for rapid spread of SARS-
396 CoV-2 variants. Our findings suggest that global SARS-CoV-2 genomic surveillance strategies and
397 capacity are diverse, and may be limited in some regions, especially in the context of the global
398 spread and dominance of variants of concern. The gap still exists in sequencing availability and
399 magnitude, therefore international efforts are needed to address some genomic bottleneck, such as
400 lack of tracking representative samples, lack of sequencing capacity, strict regulations about data
401 sharing, and lack of funding^{46,47}.

402 **Contributors**

403 H.Y. designed and supervised the study. Z.C., J.Z., Y.T., R.S, X.X., Y.W., W.L., S.G., and Z.Z.
404 collected data. Z.C., J.Z., Y.T., R.S, and X.X. prepared the tables. Z.C. analysed data, prepared the
405 figures, and wrote the first draft of the manuscript. A.S.A., X.C., J.Y., D.T.L., D.B.D, and H.Y.
406 interpreted the results and revised the content critically. All authors contributed to review and revision
407 and approved the final manuscript as submitted and agree to be accountable for all aspects of the work.
408

409 **Declaration of interests**

410 H.Y. has received research funding from Sanofi Pasteur, GlaxoSmithKline, Yichang HEC Changjiang
411 Pharmaceutical Company, and Shanghai Roche Pharmaceutical Company. None of those research
412 funding is related to COVID-19. All other authors report no competing interests.
413

414 **Role of the funding source**

415 The funder had no role in study design, data collection, data analysis, data interpretation, or writing of
416 the report. The corresponding author had full access to all the data in the study and had final
417 responsibility for the decision to submit for publication.
418

419 **Data sharing**

420 All datasets generated and analysed are available in the Article and Appendix. Any additional
421 information is available from the lead contact upon request.
422

423 **Acknowledgments**

424 We thank all the authors contributed to generating and sharing sequences to GISAID, GenBank,
425 National Genomics Data Center, National Microbiology Data Center, and China National GeneBank.
426 We thank all the authors contributed to generating and sharing aggregated data to TESSy in ECDC.
427 We thank Dr.Lina Ma, Dr.Yiming Bao, and their team from China National Center for
428 Bioinformation for building 2019nCoV dataset. We thank Junbo Chen, Xiaowei Deng from Fudan
429 University. This study was funded by Key Program of the National Natural Science Foundation of
430 China (grant no. 82130093 to HJY) and the US National Institutes of Health (R01 AI135115 to DTL
431 and ASA; KL2TR001448 to DD).

432 **Figure legends**

433 **Fig 1. Global SARS-CoV-2 genomic surveillance, sequencing availability, and publicly deposited**
434 **genomic data.**

435 (A) The global distribution of three strategies of SARS-CoV-2 genomic surveillance. (B) The global
436 availability of SARS-CoV-2 sequencing, countries with a high level of availability represent the
437 ability to perform in-country SARS-CoV-2 sequencing alone. (C) The weekly number of publicly
438 deposited SARS-CoV-2 genomic data by region. (D) Cumulative number of publicly deposited
439 SARS-CoV-2 genomic data by countries as of 20 August 2021. (E) The weekly proportion of
440 infections sequenced by region. (F) Cumulative proportion of infections sequenced by countries as of
441 20 August 2021, defined as the proportion of cumulative isolates sequences to the cumulative
442 confirmed cases. The number of sequences for the most recent weeks might be incomplete due to a
443 time delay between collecting specimens and sequencing submission. Data unavailable, includes
444 those locations that not belonged to 194 Member States or had not applicable data.

445

446 **Fig 2. The public availability extent of SARS-CoV-2 genomic data to publicly repositories by**
447 **country and the variants of concern.**

448 The public availability extent was defined as the ratio of the cumulative number of variants in
449 publicly repositories to the official reported number of variants within the same period. The ratio of
450 exceed 100% in some countries might be due to the delay in the official report of the sequencing
451 results or the incomplete official reporting system. In view of the availability of official data, the
452 cumulative number of variants in different countries corresponds to different time periods, with
453 detailed information in Table S6. Sequence without date of specimen collection in publicly
454 repositories is not included in our analysis. The variant data for China only includes cases reported by
455 mainland China. The officially reported number of Alpha variants might contain those confirmed by
456 the PCR screen assay. The values beneath the country names show numbers of cumulative variant as
457 of one specific week: variants in publicly repositories/official reported variants.

458

459 **Fig 3. The earliest identification of Alpha, Beta, Gamma, and Delta variant in each country.**

460 If information about the earliest sampling date was unavailable but that of the earliest reporting date
461 was available, we extrapolated the sampling date by using a fixed three-week lag from sample
462 collection to reporting. Countries with darker red indicate earlier samples, and with darker blue refers
463 to later samples. The white areas represent countries with variants unreported or data unavailable.

464

465 **Fig 4. The prevalence and temporal dynamics of reference strains and four SARS-CoV-2 VOCs**
466 **by country.**

467 The date presented in the top refers to the range of date of specimen collection. The prevalence was
468 defined as the proportion of the strain number (reference strains or variants) to the total number of
469 sequences generated in the same unit of time. Reference strains includes lineage A, A.1, B, and B.1;
470 the sub-lineages of four VOCs are aggregated with the parent lineages. The grey areas represent
471 countries with no COVID-19 epidemic, or no performing sequencing, or no uploading genomic data
472 to publicly repositories.

473

474 **Fig 5. The number and proportion of SARS-CoV-2 variants by region and time.**

475 The line and point in the left figure correspond to the y-axis on the right. The sub-lineages of four
476 VOCs are aggregated with the parent lineages; designated Variants of Interest (VOIs) included
477 lineage B.1.525, B.1.526, B.1.617.1, C.37, B.1.621 and their sub-lineages; other lineages included
478 reference strains and other variants. Data used here are derived from the publicly repositories and
479 aggregated dataset, with priority given to the one with highest sequenced numbers in a specific week.

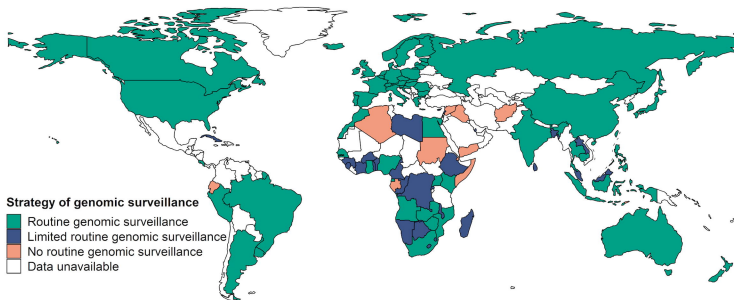
480 Reference

- 481 1. Galloway SE, Paul P, MacCannell DR, et al. Emergence of SARS-CoV-2 B.1.1.7 Lineage -
482 United States, December 29, 2020-January 12, 2021. *MMWR Morbidity and mortality weekly report*
483 2021; **70**(3): 95-9.
- 484 2. Davies NG, Abbott S, Barnard RC, et al. Estimated transmissibility and impact of SARS-
485 CoV-2 lineage B.1.1.7 in England. *Science (New York, NY)* 2021; **372**(6538).
- 486 3. Volz E, Mishra S, Chand M, et al. Assessing transmissibility of SARS-CoV-2 lineage B.1.1.7
487 in England. *Nature* 2021; **593**(7858): 266-9.
- 488 4. Wall EC, Wu M, Harvey R, et al. AZD1222-induced neutralising antibody activity against
489 SARS-CoV-2 Delta VOC. *Lancet (London, England)* 2021; **398**(10296): 207-9.
- 490 5. World Health Organization. Tracking SARS-CoV-2 variants. 2021.
491 <https://www.who.int/en/activities/tracking-SARS-CoV-2-variants/> (accessed 18 Jul 2021).
- 492 6. Campbell F, Archer B, Laurenson-Schafer H, et al. Increased transmissibility and global
493 spread of SARS-CoV-2 variants of concern as at June 2021. *Euro surveillance : bulletin European sur*
494 *les maladies transmissibles = European communicable disease bulletin* 2021; **26**(24).
- 495 7. Ito K, Piantham C, Nishiura H. Predicted dominance of variant Delta of SARS-CoV-2 before
496 Tokyo Olympic Games, Japan, July 2021. *Euro surveillance : bulletin European sur les maladies*
497 *transmissibles = European communicable disease bulletin* 2021; **26**(27).
- 498 8. Scientific Advisory Group for Emergencies (SAGE). SPI-M-O: Consensus Statement on
499 COVID-19. 2021.
500 [https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/986](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/986709/S1237_SPI-M-O_Consensus_Statement.pdf)
501 [709/S1237_SPI-M-O_Consensus_Statement.pdf](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/986709/S1237_SPI-M-O_Consensus_Statement.pdf).
- 502 9. Dhar MS, Marwal R, Radhakrishnan VS, et al. Genomic characterization and Epidemiology
503 of an emerging SARS-CoV-2 variant in Delhi, India. *medRxiv* 2021: 2021.06.02.21258076.
- 504 10. Chen X, Chen Z, Azman AS, et al. Neutralizing antibodies against SARS-CoV-2 variants
505 induced by natural infection or vaccination: a systematic review and pooled meta-analysis. *Clinical*
506 *infectious diseases : an official publication of the Infectious Diseases Society of America* 2021.
- 507 11. Abu-Raddad LJ, Chemaitelly H, Butt AA. Effectiveness of the BNT162b2 Covid-19 Vaccine
508 against the B.1.1.7 and B.1.351 Variants. *The New England journal of medicine* 2021; **385**(2): 187-9.
- 509 12. Sheikh A, McMenemy J, Taylor B, Robertson C. SARS-CoV-2 Delta VOC in Scotland:
510 demographics, risk of hospital admission, and vaccine effectiveness. *The Lancet* 2021; **397**(10293):
511 2461-2.
- 512 13. Li XN, Huang Y, Wang W, et al. Efficacy of inactivated SARS-CoV-2 vaccines against the
513 Delta variant infection in Guangzhou: A test-negative case-control real-world study. *Emerging*
514 *microbes & infections* 2021: 1-32.
- 515 14. Vogels CBF, Breban MI, Ott IM, et al. Multiplex qPCR discriminates variants of concern to
516 enhance global surveillance of SARS-CoV-2. *PLoS biology* 2021; **19**(5): e3001236.
- 517 15. Wu F, Zhao S, Yu B, et al. A new coronavirus associated with human respiratory disease in
518 China. *Nature* 2020; **579**(7798): 265-9.
- 519 16. Hodcroft EB, Zuber M, Nadeau S, et al. Spread of a SARS-CoV-2 variant through Europe in
520 the summer of 2020. *Nature* 2021.
- 521 17. Elbe S, Buckland-Merrett G. Data, disease and diplomacy: GISAID's innovative contribution
522 to global health. *Global challenges (Hoboken, NJ)* 2017; **1**(1): 33-46.
- 523 18. European Centre for Disease Prevention and Control. Sequencing of SARS-CoV-2: first
524 update, 2021.
- 525 19. Bugembe DL, Phan MVT, Ssewanyana I, et al. Emergence and spread of a SARS-CoV-2
526 lineage A variant (A.23.1) with altered spike protein in Uganda. *Nature microbiology* 2021.
- 527 20. Hodcroft EB, Zuber M, Nadeau S, et al. Spread of a SARS-CoV-2 variant through Europe in
528 the summer of 2020. *Nature* 2021.
- 529 21. Wilkinson E, Giovanetti M, Tegally H, et al. A year of genomic surveillance reveals how the
530 SARS-CoV-2 pandemic unfolded in Africa. *medRxiv* 2021: 2021.05.12.21257080.
- 531 22. Gong Z, Zhu JW, Li CP, et al. An online coronavirus analysis platform from the National
532 Genomics Data Center. *Zoological research* 2020; **41**(6): 705-8.

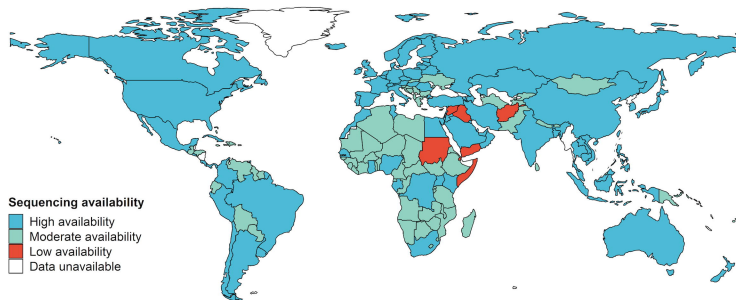
- 533 23. Shu Y, McCauley J. GISAID: Global initiative on sharing all influenza data - from vision to
534 reality. *Euro surveillance : bulletin Europeen sur les maladies transmissibles = European*
535 *communicable disease bulletin* 2017; **22**(13).
- 536 24. Hatcher EL, Zhdanov SA, Bao Y, et al. Virus Variation Resource - improved response to
537 emergent viral outbreaks. *Nucleic acids research* 2017; **45**(D1): D482-d90.
- 538 25. National Genomics Data Center Members and Partners. Database Resources of the National
539 Genomics Data Center in 2020. *Nucleic acids research* 2020; **48**(D1): D24-d33.
- 540 26. Shi W, Qi H, Sun Q, et al. gcMeta: a Global Catalogue of Metagenomics platform to support
541 the archiving, standardization and analysis of microbiome data. *Nucleic acids research* 2019; **47**(D1):
542 D637-d48.
- 543 27. Xiao SZ, Armit C, Edmunds S, et al. Increased interactivity and improvements to the
544 GigaScience database, GigaDB. *Database : the journal of biological databases and curation* 2019;
545 **2019**.
- 546 28. Song S, Ma L, Zou D, et al. The Global Landscape of SARS-CoV-2 Genomes, Variants, and
547 Haplotypes in 2019nCoV. *Genomics, proteomics & bioinformatics* 2020; **18**(6): 749-59.
- 548 29. World Health Organization. WHO Coronavirus (COVID-19) Dashboard.
549 <https://covid19.who.int/info/>.
- 550 30. World Health Organization. Operational considerations to expedite genomic sequencing
551 component of GISRS surveillance of SARS-CoV-2, 2021.
- 552 31. Africa CDC. Africa Pathogen Genomics Initiative (PGI). 2021.
553 <https://africacdc.org/institutes/africa-pathogen-genomics-initiative/> (accessed Jul 18 2021).
- 554 32. Paul P, France AM, Aoki Y, et al. Genomic Surveillance for SARS-CoV-2 Variants
555 Circulating in the United States, December 2020-May 2021. *MMWR Morbidity and mortality weekly*
556 *report* 2021; **70**(23): 846-50.
- 557 33. Brown KA, Gubbay J, Hopkins J, et al. S-Gene Target Failure as a Marker of Variant B.1.1.7
558 Among SARS-CoV-2 Isolates in the Greater Toronto Area, December 2020 to March 2021. *Jama*
559 *2021*; **325**(20): 2115-6.
- 560 34. Crawford DC, Williams SM. Global variation in sequencing impedes SARS-CoV-2
561 surveillance. *PLoS genetics* 2021; **17**(7): e1009620.
- 562 35. World Health Organization. Genomic sequencing of SARS-CoV-2: A guide to
563 implementation for maximum impact on public health, 2021.
- 564 36. World Health Organization. Guidance for surveillance of SARS-CoV-2 variants: Interim
565 guidance. 2021. <https://www.who.int/publications/i/item/WHO-2019-nCoV-surveillance-variants>
566 (accessed August 26 2021).
- 567 37. Page AJ, Mather AE, Le-Viet T, et al. Large-scale sequencing of SARS-CoV-2 genomes
568 from one region allows detailed epidemiology and enables local outbreak management. *Microbial*
569 *genomics* 2021; **7**(6).
- 570 38. Black A, MacCannell DR, Sibley TR, Bedford T. Ten recommendations for supporting open
571 pathogen genomic analysis in public health. *Nature medicine* 2020; **26**(6): 832-41.
- 572 39. Capoferri AA, Shao W, Spindler J, Coffin JM, Rausch JW, Kearney MF. 2020 SARS-CoV-2
573 diversification in the United States: Establishing a pre-vaccination baseline. *medRxiv* 2021:
574 2021.06.01.21258185.
- 575 40. Robishaw JD, Alter SM, Solano JJ, et al. Genomic surveillance to combat COVID-19:
576 challenges and opportunities. *The Lancet Microbe* 2021.
- 577 41. Li B, Deng A, Li K, et al. Viral infection and transmission in a large, well-traced outbreak
578 caused by the SARS-CoV-2 Delta variant. *medRxiv* 2021: 2021.07.07.21260122.
- 579 42. Meng Z, Jianpeng X, Aiping D, et al. Transmission Dynamics of an Outbreak of the COVID-
580 19 Delta Variant B.1.617.2 — Guangdong Province, China, May–June 2021. *China CDC Weekly*
581 *2021*; **3**(27): 584-6.
- 582 43. Pung R, Mak TM, Kucharski AJ, Lee VJ. Serial intervals in SARS-CoV-2 B.1.617.2 variant
583 cases. *The Lancet*.
- 584 44. Brown CM, Vostok J, Johnson H, et al. Outbreak of SARS-CoV-2 Infections, Including
585 COVID-19 Vaccine Breakthrough Infections, Associated with Large Public Gatherings - Barnstable
586 County, Massachusetts, July 2021. *MMWR Morbidity and mortality weekly report* 2021; **70**(31):
587 1059-62.

- 588 45. Krause PR, Fleming TR, Longini IM, et al. SARS-CoV-2 Variants and Vaccines. *New*
589 *England Journal of Medicine* 2021; **385**(2): 179-86.
- 590 46. Hodcroft EB, De Maio N, Lanfear R, et al. Want to track pandemic variants faster? Fix the
591 bioinformatics bottleneck. *Nature* 2021; **591**(7848): 30-3.
- 592 47. Wadman M. United States rushes to fill void in viral sequencing. 2021; **371**(6530): 657-8.
593

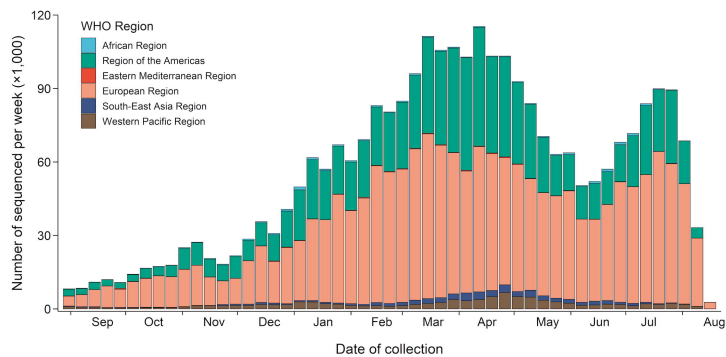
A



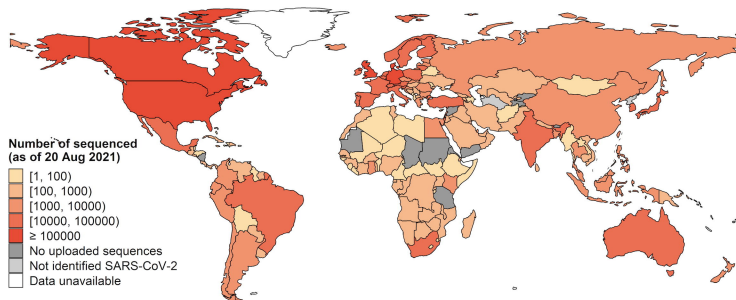
B



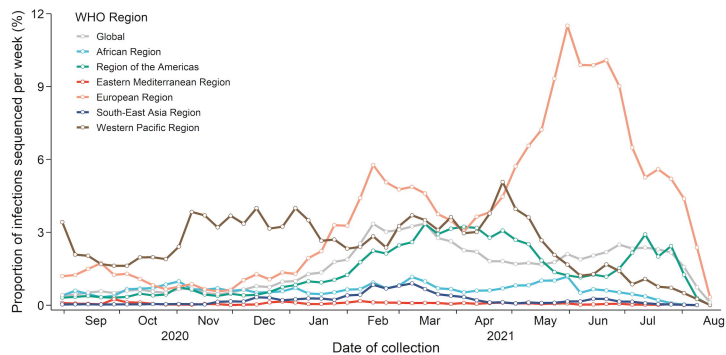
C



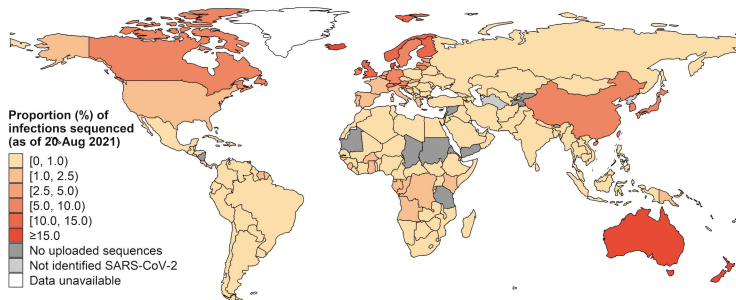
D

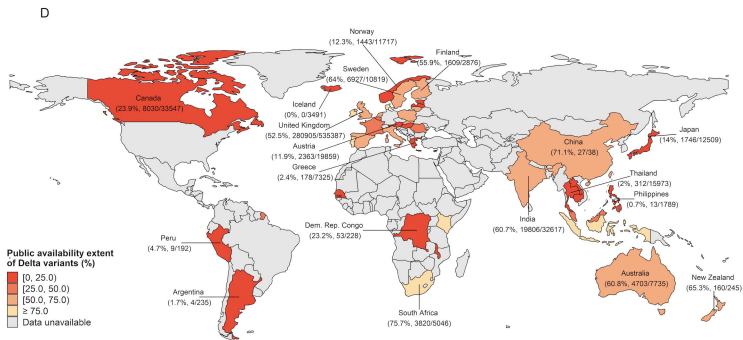
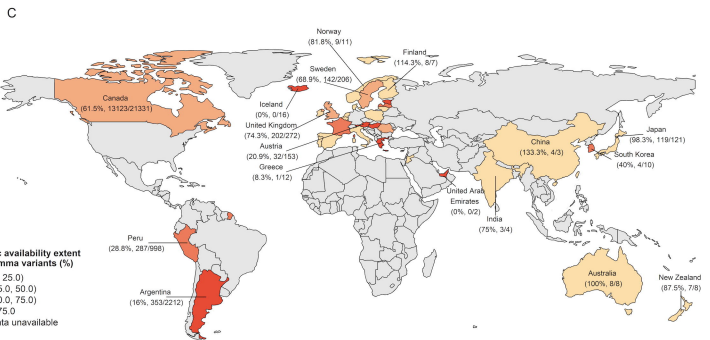
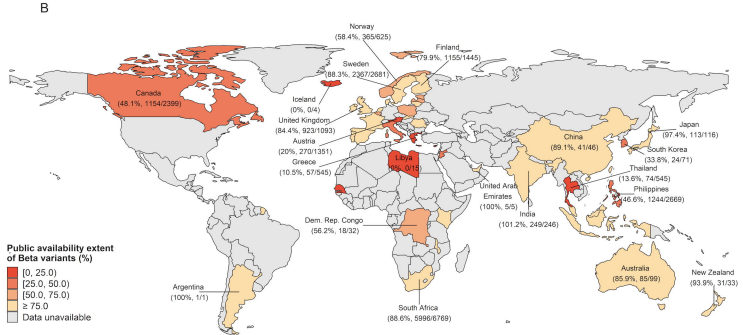
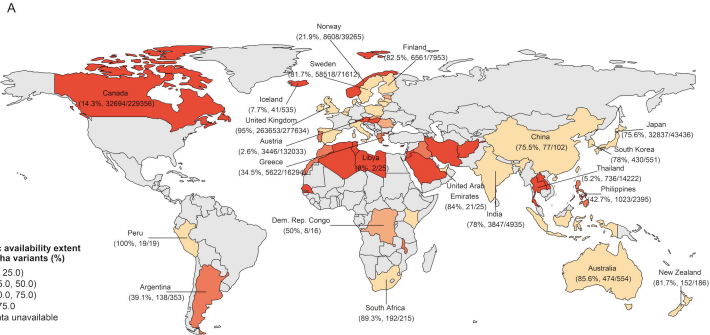


E

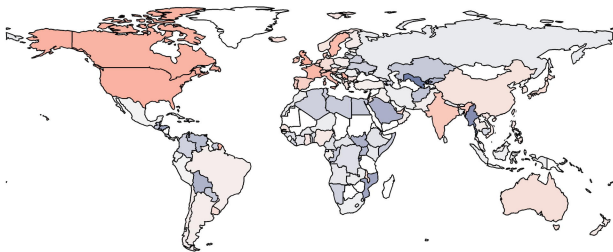


F

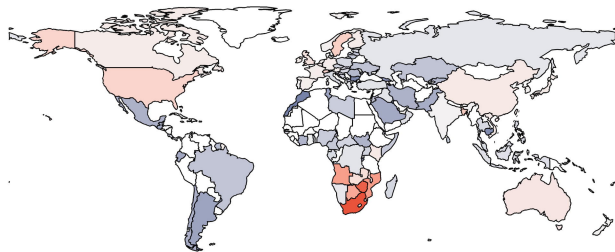




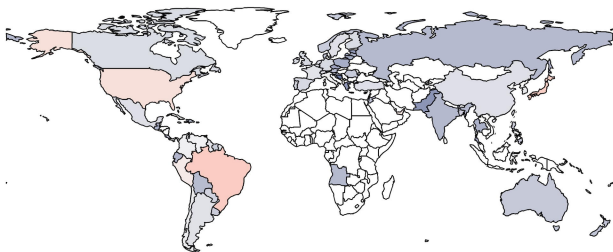
A. Alpha variant



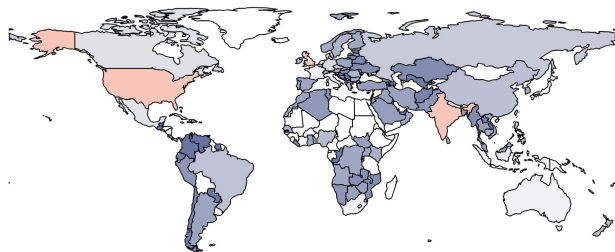
B. Beta variant



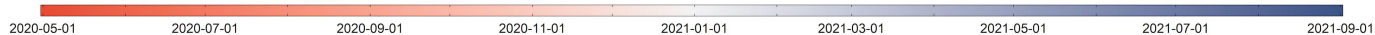
C. Gamma variant



D. Delta variant



Earliest sampling date of VOCs



Jan-Dec, 2020

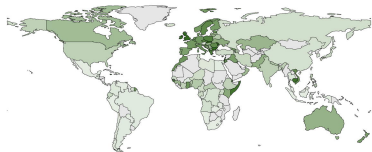
Jan-Mar, 2021

Apr-Jun, 2021

Jul-Aug, 2021



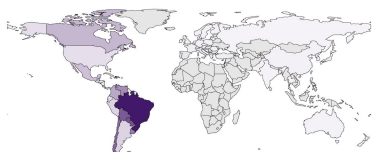
Prevalence of reference strains
0 25 50 75 100



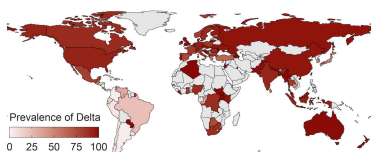
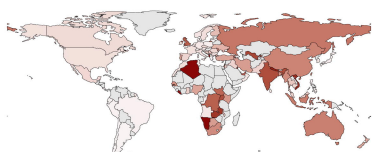
Prevalence of Alpha
0 25 50 75 100



Prevalence of Beta
0 25 50 75 100



Prevalence of Gamma
0 25 50 75 100



Prevalence of Delta
0 25 50 75 100

