

Global landscape of SARS-CoV-2 genomic surveillance, public availability extent of genomic data, and epidemic shaped by variants

Hongjie Yu (✉ yhj@fudan.edu.cn)

Fudan University <https://orcid.org/0000-0002-6335-5648>

Zhiyuan Chen

Fudan University

Andrew Azman

Department of Epidemiology, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA

Xinhua Chen

Fudan University

Junyi Zou

Fudan University

Yuyang Tian

Fudan University

Ruijia Sun

Fudan University

Xiangyanyu Xu

Fudan University

Yani Wu

Fudan University

Wanying Lu

Fudan University

Shijia Ge

Fudan University

Zeyao Zhao

Fudan University

Juan Yang

Fudan University

Daniel Leung

University of Utah School of Medicine, Salt Lake City, UT, USA

Daryl Domman

Center for Global Health, Department of Internal Medicine, University of New Mexico Health Sciences Center, New Mexico, USA

Article

Keywords: SARS-CoV-2, genomic data, variants, genomic surveillance

Posted Date: September 29th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-927070/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Abstract

Genomic surveillance has shaped our understanding of SARS-CoV-2 variants, which have proliferated globally in 2021. We collected country-specific data on SARS-CoV-2 genomic surveillance, sequencing capabilities, public genomic data from multiple public repositories, and aggregated publicly available variant data. Then, different proxies were used to estimate the sequencing coverage and public availability extent of genomic data, in addition to describing the global dissemination of variants. We found that the COVID-19 global epidemic clearly featured increasing circulation of Alpha since the start of 2021, which was rapidly replaced by the Delta variant starting around May 2021. SARS-CoV-2 genomic surveillance and sequencing availability varied markedly across countries, with 63 countries performing routine genomic surveillance and 79 countries with high availability of SARS-CoV-2 sequencing. We also observed a marked heterogeneity of sequenced coverage across regions and countries. Across different variants, 21-46% of countries with explicit reporting on variants shared less than half of their variant sequences in public repositories. Our findings indicated an urgent need to expand sequencing capacity of virus isolates, enhance the sharing of sequences, the standardization of metadata files, and supportive networks for countries with no sequencing capability.

Background

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

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

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

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

Results

We classified genomic surveillance strategies for 105 countries, including 72.3% (34/47) of WHO-defined Africa Region countries, 58.5% (31/53) of European Region countries, 61.9% (13/21) of countries in the Eastern Mediterranean Region, 31.4% (11/35) of countries in Americas Region, 37.0% (10/27) of countries in the Western Pacific Region, and 54.5% (6/11) of countries in the South East Asia Region (Table S4). We downloaded a total of 3.8 million deduplicated SARS-CoV-2 sequence samples from publicly repositories corresponding to samples collected between 1 December 2019 to 15 September 2021 in 164 countries. Additionally, we collected official aggregated data of variants from 55 countries, and extra data for the first identification of COVID-19 variants from 33 countries.

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

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

with data accessible, the availability of SARS-CoV-2 sequences was high in 79 countries, moderate in 80 countries, and low in 8 countries (Fig. 1B).

Sequencing coverage of SARS-CoV-2 confirmed cases

The European (57.3%) and Americas (35.3%) Regions uploaded the most SARS-CoV-2 sequences to public repositories, with marked intra-region heterogeneity across countries, with a range from 3 (Saint Kitts and Nevis) to 1.1 million (United States) as of 15 September 2021 (Fig. 1D).

Since September 2020, no more than 4.0% of global confirmed SARS-CoV-2 infections were sequenced, with a relatively high proportion of infections sequenced in March (3.4%) and early August (3.8%). In any week, the Africa, South East Asia and Eastern Mediterranean Regions (Fig. 1E) sequenced no more than 1.5% of cases. At the end of May 2021 (start of Delta widely spreading), Europe had the highest sequenced proportion of 9.6%, followed by Western Pacific (2.1%), Americas (1.4%), Africa (1.1%), South-East Asia (0.1%), and Eastern Mediterranean (0.09%) (Fig. 1E). At the country-level, higher rates of sequencing were observed in Iceland, New Zealand, Denmark, Australia, Luxembourg, Norway, Finland, and United Kingdom, all of which had at least 10% reported infections sequenced as of 15 September 2021. In addition, almost all countries in the Africa and Eastern Mediterranean had sequenced less than 2.5% of confirmed infections, except for Gambia (6.3%) and Djibouti (3.1%) (Fig. 1F).

Public availability extent of SARS-CoV-2 variant sequences

The public availability extent of SARS-CoV-2 genomic data varied across variants and countries. Overall, among countries with aggregated data on the number of variant infections, less than half of sequences of Alpha, Beta, Gamma, and Delta were publicly available in 38.5% (20/52), 20.9% (9/43), 25.8% (8/31), and 45.5% (20/44) countries, respectively (Fig. 2). However, the result for Alpha might be influenced by SGTF detected via PCR. The public availability extent of Delta variants across countries ranged from 0.0% (Cyprus, Hungary, Laos) to 100.0% (New Zealand, Indonesia, Netherlands, South Africa). At the country level, low sharing proportions (less than 50%) across all VOCs was found in several countries, including Austria, Cyprus, Greece, Hungary, Philippines, Thailand and Senegal. For example, the sharing proportion of Alpha, Beta, and Delta in Thailand was 10.5%, 15.0%, and 4.0%, respectively, which indicated more than 85.0% genomic data of variants were not uploaded.

Moreover, incomplete metadata attached to GISAID sequences was common globally, with about two thirds of sequences missing demographic information (age and sex), and more than 95% of that missing clinical information (e.g., symptom history, clinical outcome, and vaccination status) (Table S5). High-income regions tend to have low information completeness, especially in the European Region, where less than 25.0% and 3.0% of sequences had demographic and clinical information, respectively, significantly lower than other Regions ($P < 0.0001$, Chi-square test). Besides, 93.4% sequences were reported by subnational geographies, with a low proportion in Western Pacific (58.2%) and Eastern Mediterranean (76.8%).

Earliest identification of SARS-CoV-2 variants across regions

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

Global and regional spread of SARS-CoV-2 variants

The number of new VOC cases dramatically increased until April 2021, with a peak weekly value of about 100,000 VOC cases sequenced in which most of them were Alpha variants (Fig. 5A). Subsequently, another peak of weekly new VOC case occurred in early August 2021, but with a large amount of Delta variants. The number of VOC cases may be an underestimate for the most recent weeks due to a collection-to-report time delay. Notably, this increase was also accompanied by the increase in the volume of new sequenced cases and new COVID-19 confirmed cases.

The global prevalence of reference (non-variant) strains fell into a low level of 0.6% in the period of Jul-Sep 2021, compared with 14.0% of that in 2020 (Fig. 4). Globally, the COVID-19 pandemic was driven by the circulation of Alpha at the start of 2021, with an average prevalence of 50.4% in the first quarter of 2021. Alpha variants continued to outcompete other strains in the second quarter of 2021, accounting for 59.8% of the contemporary lineages (Fig. 4). However, the rapid global rise of the Delta variant began in May 2021, reaching a global prevalence of nearly 98.8% at the end of August 2021 (Fig. 5B). In contrast, Beta and Gamma variants remained at low prevalence (Fig. 4). Additionally, the shifting of predominant variants from Alpha to Delta first occurred in South-East Asia where the proportion of Delta exceeded 60.0% in April 2021 (Fig. 5J).

Discussion

Our study characterized the global diversity of genomic surveillance strategies and sequencing availability, sequenced coverage of SARS-CoV-2 cases, public availability extent of variant sequences, as well as current epidemic situation of SARS-CoV-2 variants. We found that genomic surveillance strategies were globally heterogenous, with limited or no routine surveillance among many countries in the Africa and Eastern Mediterranean Regions. Our analysis of publicly deposited SARS-CoV-2 sequences implied that the sequenced coverage is low in most countries, with a low proportion of VOCs sequences shared to public repositories. The pervasive spread of Alpha and Delta variants further highlights the threat of SARS-CoV-2 mutations despite the availability of vaccines in many countries.

The diversity of SARS-CoV-2 genomic surveillance between countries is associated with country-specific priorities (e.g., surveillance objectives, targeted monitoring, or event-/risk-based sequencing) and available resources. ECDC recommends population-based and/or targeted sampling strategies (e.g., imported cases, cluster cases, and potential vaccine escapers) for genomic surveillance, which could provide a more representative estimate of the relative prevalence of variants. Notably, several countries, many of which are classified as low- or lower middle-income countries by the World Bank, lack genomic surveillance data, likely due to limitations in infrastructure capacity and resources. However, even some countries classified as high-income, have suffered from a slow and inconsistent process of adopting genomics-based surveillance²². Establishment of reference laboratories and networks to provide sequencing services for countries without established sequencing capacity may enable improved detection and tracking of emerging variants worldwide.

The detection of most variants relies on the full-length or partial genomic sequencing, but the sequences only become available for the global community when the laboratories have established sequencing capacity, willing to share, and legally allowed to upload them. The discrepancies in sharing was observed in each region, which confirmed that some countries are sequencing but are not uploading. However, our study observed a sharing extent of exceed 100% exists in some countries, likely due to delays in the official reporting of sequencing results, or the incomplete official reporting system. The timely sharing of those enables to adequately contextualize local data when looking at introductions and examine transmission routes, as well as to look for sites of repeated mutations that can guide laboratory work on characterizing those mutations effects on therapeutics and vaccine efficacy. The underlying reasons why some countries didn't share might be related with the distrust for publicly repositories in the concept of data security.

We found relatively low completeness of demographic and clinical characteristics in metadata accompanying uploaded sequences. Our analysis revealed that high-income countries frequently did not share demographic information. A possible reason for this is that these regions may having more restrictive data privacy/laws preventing/discouraging the release of this information. Genomic data coupled with those additional data can maximize the utility of genomic sequences in rapid scientific discovery during this pandemic, which are valuable for in-depth epidemiological analyses to characterize risk factors, clinical severity, and other public health risk of variants²³⁻²⁵. Therefore, it's vital to optimize the sharing of information in a secure and trusted channel in the context of protect patient anonymity and in accordance of local regulations²⁶. In addition, decreasing the lag between sample collection to deposition of these sequences²⁷, including the timely sharing and standardizing of metadata^{18,23,28}, may facilitate the design and development of treatment and prevention strategies by policy-makers²⁹.

An important role of genomic surveillance is to investigate the spread and dynamics of SARS-CoV-2 variants. Amidst the emergence of different variants, the current dominance of the Delta variant suggests that it may possess higher fitness than other variants, which might be associated with a combination of higher virus load, and shorter incubation period and serial intervals³⁰⁻³². The decrease in real-world

vaccine effectiveness against Beta variants¹¹ and increasing breakthrough cases with Delta variants³³ underlines the importance of determining the local or regional patterns of variant spread, including the need to develop new or modified vaccines to achieve adequate protection³⁴.

Our results should be interpreted in view of several limitations. First, the lack of data from some countries limited our global mapping. The data completeness and quality could be impacted by key steps in the surveillance or reporting, including differences in diagnostic criteria, under-reporting, delayed reporting, and reporting methods. The inconsistent diagnostic criteria of variants might cause sampling bias, especially when adopting PCR assay to detect Alpha variant owing to its non-specificity³⁵. We did an extensive search to collect multi-source data and chose the aggregated data with a priority to sequencing results rather than PCR-screening results. Second, the analysis of global and national spread could be biased as data from public repositories or aggregated dataset are not always representative of the variants circulating in the regions, especially for the regions with relatively limited sequencing capacity or with investigating outbreak-based events. Therefore, the global prevalence of variants may be biased due to the uneven sequencing across the regions. Indeed, it's difficult to obtain a truly representative and random sample, and how to understand these biases will become important²⁴. Lastly, the detailed demographical, epidemiological and clinical information about variant cases cannot be accessed, which limited our further epidemiological analysis about variant spread.

In conclusion, our study provides a landscape for genomic surveillance, the global coverage of sequencing and public availability of sequences, as well as the evidence for rapid spread of SARS-CoV-2 variants. Our findings suggest that global SARS-CoV-2 genomic surveillance strategies and capacity are diverse, and may be limited in some regions, especially in the context of the global spread and dominance of variants of concern. The gap still exists in sequencing availability and magnitude, therefore international efforts are needed to address some genomic bottleneck, and more work are needed to be done in defining the ideal sampling schemes for different purposes and sharing these data in public repositories to allow for further rapid scientific discovery.

Methods

Data sources and collection

Through extracting country-specific data from multiple publicly available sources, we built three datasets of genomic surveillance, deposited genomic data in publicly repositories, and official aggregated genomic data as of 15 September 2021.

Dataset of genomic surveillance

Each country's genomic surveillance strategy and sequencing availability was gathered from searches of the websites of regional WHO, the country's Ministry of Health, CDC, local academic partners, and official news, supplemented by a literature search (appendix). Data extracted included the overall surveillance

strategy, sequencing availability, target population, sampling method, diagnostic criteria, and sequenced volume. Given that the surveillance strategy and density may change with time, we only gathered information on the most recent surveillance strategy (as of 15 September 2021).

Dataset of SARS-CoV-2 sequences in publicly repositories

SARS-CoV-2 sequences along with related metadata file were downloaded from an online coronavirus analysis platform from the National Genomics Data Center (NGDC)³⁶, where has merged and deduplicated sequences that deposited in GISAID³⁷, GenBank³⁸, National Genomics Data Center³⁹, National Microbiology Data Center⁴⁰, and China National GeneBank⁴¹. Detailed process of integration and deduplication was previously reported⁴². An acknowledgement table for those contributing to this work is available in the Appendix.

Official aggregated dataset

To gain additional insights regarding the public availability (sharing) extent of genomic data of SARS-CoV-2 variants, we extracted country-specific, variant-specific, and time-specific aggregated data on the number of SARS-CoV-2 variant cases from official websites, using the same sources as above, except for the literature source. The search was done by either directly locating to the official website for each country or indirectly searching in search engines (Google, Bing, Baidu) by using the terms “variant” and country name. To supplement the aggregated data of variants that we collected, we also downloaded the aggregated data with a valid denominator (namely, the number of isolates sequenced is reasonable) from the European Surveillance System (TESSy). The aggregated dataset included country name, date of report or collection, new or cumulative numbers of different SARS-CoV-2 variant cases and total sequenced cases, as well as the new number of confirmed cases. COVID-19 epidemic data were derived from WHO⁴³. Considering that the diagnostic criteria of SARS-CoV-2 variants varies in different countries, the general principle for collecting aggregated data was to give priority to the results based on whole and partial genome sequencing instead of those based on a PCR assay.

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

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

Data analysis

We used the variant naming system proposed by WHO, where four VOCs and five Variants of Interest (VOIs, included Eta, Iota, Kappa, Lambda, and Mu) had been designated as of 20 September 2021⁵. Our study focused on analyses on four VOCs in the 194 Member States of WHO. We did not integrate data

from the overseas territories into that country's data. Given most countries weekly released aggregated data of variant, most of our analyses were performed on a weekly basis.

To characterize the global diversity of SARS-CoV-2 genomic surveillance, we classified the surveillance strategy of each country into three categories: 1) routine genomic surveillance; 2) limited routine genomic surveillance; and 3) no routine genomic surveillance. Routine genomic surveillance was defined as conducting nationwide genomic sequencing, coupled with at least 150 specimens per week or 10% of all samples sequenced²⁸. The classification of most African countries was from a pre-defined version from African CDC⁴⁴, while other countries were subsequently defined according to the definition criteria (Table S3). As we could not identify information on the surveillance strategy for some countries through public sources, we also classified three extra categories according to the public availability/ability of genomic sequencing: 1) high availability; 2) moderate availability; and 3) low availability (Table S3).

The process of data cleaning in publicly repositories and aggregated dataset is shown in the Appendix. For data that was reported in aggregate, when the date of sample collection was not available, we assumed a fixed three-week lag⁴⁵ from sample collection to reporting unless other country-specific information was available to inform this extrapolation.

The sequencing coverage was inferred by using the percent of cumulative positives sequenced as a proxy, defined by the ratio of the number of isolates sequenced to the number of confirmed cases in the same unit of time. Given that not all sequences will be uploaded to genomic repositories, we analysed the public availability extent of genomic data by comparing the cumulative number of variants between public repositories and official aggregated datasets. Since the Alpha variant had a characteristic deletion of amino acids 69–70 in the SGTF that can be detected via a widely used PCR assay⁴⁶, we performed this analysis across VOCs.

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

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

The comparison between ratios of two groups was done using the Chi-square test for, and differences were considered statistically significant at P -value < 0.05 . All statistical analyses and visualizations were done using R (version 4.0.2).

Declarations

Contributors

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. collected data. Z.C., J.Z., Y.T., R.S, and X.X. prepared the tables. Z.C. analysed data, prepared the 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. interpreted the results and revised the content critically. All authors contributed to review and revision and approved the final manuscript as submitted and agree to be accountable for all aspects of the work.

Declaration of interests

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

Role of the funding source

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

Data sharing

Any additional information is available from the lead contact upon reasonable request.

Acknowledgments

We thank all the authors contributed to generating and sharing sequences to GISAID, GenBank, National Genomics Data Center, National Microbiology Data Center, and China National GeneBank. We thank all the authors contributed to generating and sharing aggregated data to TESSy in ECDC. We thank Dr.Lina

Ma, Dr.Yiming Bao, and their team from China National Center for Bioinformation for building 2019nCoV dataset. We thank Junbo Chen, Xiaowei Deng, and Xuemei Yan from Fudan University. This study was funded by Key Program of the National Natural Science Foundation of China (grant no. 82130093 to HJY) and the US National Institutes of Health (R01 AI135115 to DTL and ASA; KL2TR001448 to DD).

References

1. Galloway, S.E. *et al.* Emergence of SARS-CoV-2 B.1.1.7 Lineage - United States, December 29, 2020-January 12, 2021. *MMWR Morb Mortal Wkly Rep* **70**, 95–99 (2021).
2. Davies, N.G. *et al.* Estimated transmissibility and impact of SARS-CoV-2 lineage B.1.1.7 in England. *Science* **372**(2021).
3. Volz, E. *et al.* Assessing transmissibility of SARS-CoV-2 lineage B.1.1.7 in England. *Nature* **593**, 266–269 (2021).
4. Wall, E.C. *et al.* AZD1222-induced neutralising antibody activity against SARS-CoV-2 Delta VOC. *Lancet* **398**, 207–9 (2021).
5. World Health Organization. Tracking SARS-CoV-2 variants. Vol. 2021 (2021).
6. Campbell, F. *et al.* Increased transmissibility and global spread of SARS-CoV-2 variants of concern as at June 2021. *Euro Surveill* **26**(2021).
7. Ito, K., Piantham, C. & Nishiura, H. Predicted dominance of variant Delta of SARS-CoV-2 before Tokyo Olympic Games, Japan, July 2021. *Euro Surveill* **26**(2021).
8. Scientific Advisory Group for Emergencies (SAGE). SPI-M-O: Consensus Statement on COVID-19. (2021).
9. Dhar, M.S. *et al.* Genomic characterization and Epidemiology of an emerging SARS-CoV-2 variant in Delhi, India. *medRxiv*, 2021.06.02.21258076 (2021).
10. Chen, X. *et al.* Neutralizing antibodies against SARS-CoV-2 variants induced by natural infection or vaccination: a systematic review and pooled meta-analysis. *Clin Infect Dis* (2021).
11. Abu-Raddad, L.J., Chemaitelly, H. & Butt, A.A. Effectiveness of the BNT162b2 Covid-19 Vaccine against the B.1.1.7 and B.1.351 Variants. *N Engl J Med* **385**, 187–189 (2021).
12. Sheikh, A., McMenamin, J., Taylor, B. & Robertson, C. SARS-CoV-2 Delta VOC in Scotland: demographics, risk of hospital admission, and vaccine effectiveness. *The Lancet* **397**, 2461–2462 (2021).
13. Li, X.N. *et al.* Efficacy of inactivated SARS-CoV-2 vaccines against the Delta variant infection in Guangzhou: A test-negative case-control real-world study. *Emerg Microbes Infect*, 1–32 (2021).
14. Vogels, C.B.F. *et al.* Multiplex qPCR discriminates variants of concern to enhance global surveillance of SARS-CoV-2. *PLoS Biol* **19**, e3001236 (2021).

15. Wu, F. *et al.* A new coronavirus associated with human respiratory disease in China. *Nature* **579**, 265–269 (2020).
16. Hodcroft, E.B. *et al.* Spread of a SARS-CoV-2 variant through Europe in the summer of 2020. *Nature* (2021).
17. Elbe, S. & Buckland-Merrett, G. Data, disease and diplomacy: GISAID's innovative contribution to global health. *Glob Chall* **1**, 33–46 (2017).
18. European Centre for Disease Prevention and Control. Sequencing of SARS-CoV-2: first update. (2021).
19. Bugembe, D.L. *et al.* Emergence and spread of a SARS-CoV-2 lineage A variant (A.23.1) with altered spike protein in Uganda. *Nat Microbiol* (2021).
20. Hodcroft, E.B. *et al.* Spread of a SARS-CoV-2 variant through Europe in the summer of 2020. *Nature* (2021).
21. Wilkinson, E. *et al.* A year of genomic surveillance reveals how the SARS-CoV-2 pandemic unfolded in Africa. *medRxiv*, 2021.05.12.21257080 (2021).
22. Crawford, D.C. & Williams, S.M. Global variation in sequencing impedes SARS-CoV-2 surveillance. *PLoS Genet* **17**, e1009620 (2021).
23. World Health Organization. Genomic sequencing of SARS-CoV-2: A guide to implementation for maximum impact on public health. Vol. 2021 (2021).
24. World Health Organization. Guidance for surveillance of SARS-CoV-2 variants: Interim guidance. Vol. 2021 (2021).
25. Page, A.J. *et al.* Large-scale sequencing of SARS-CoV-2 genomes from one region allows detailed epidemiology and enables local outbreak management. *Microb Genom* **7**(2021).
26. Black, A., MacCannell, D.R., Sibley, T.R. & Bedford, T. Ten recommendations for supporting open pathogen genomic analysis in public health. *Nat Med* **26**, 832–841 (2020).
27. Kalia, K., Saberwal, G. & Sharma, G. The lag in SARS-CoV-2 genome submissions to GISAID. *Nat Biotechnol* (2021).
28. World Health Organization. Operational considerations to expedite genomic sequencing component of GISRS surveillance of SARS-CoV-2. (2021).
29. Robishaw, J.D. *et al.* Genomic surveillance to combat COVID-19: challenges and opportunities. *Lancet Microbe* (2021).
30. Li, B. *et al.* Viral infection and transmission in a large, well-traced outbreak caused by the SARS-CoV-2 Delta variant. *medRxiv*, 2021.07.07.21260122 (2021).
31. Meng, Z. *et al.* Transmission Dynamics of an Outbreak of the COVID-19 Delta Variant B.1.617.2 – Guangdong Province, China, May–June 2021. *China CDC Weekly* **3**, 584–586 (2021).
32. Pung, R., Mak, T.M., Kucharski, A.J. & Lee, V.J. Serial intervals in SARS-CoV-2 B.1.617.2 variant cases. *The Lancet*.

33. Brown, C.M. *et al.* Outbreak of SARS-CoV-2 Infections, Including COVID-19 Vaccine Breakthrough Infections, Associated with Large Public Gatherings - Barnstable County, Massachusetts, July 2021. *MMWR Morb Mortal Wkly Rep* **70**, 1059–1062 (2021).
34. Krause, P.R. *et al.* SARS-CoV-2 Variants and Vaccines. *New England Journal of Medicine* **385**, 179–186 (2021).
35. Hodcroft, E.B. *et al.* Want to track pandemic variants faster? Fix the bioinformatics bottleneck. *Nature* **591**, 30–33 (2021).
36. Gong, Z. *et al.* An online coronavirus analysis platform from the National Genomics Data Center. *Zool Res* **41**, 705–708 (2020).
37. Shu, Y. & McCauley, J. GISAID: Global initiative on sharing all influenza data - from vision to reality. *Euro Surveill* **22**(2017).
38. Hatcher, E.L. *et al.* Virus Variation Resource - improved response to emergent viral outbreaks. *Nucleic Acids Res* **45**, D482-d490 (2017).
39. National Genomics Data Center Members and Partners. Database Resources of the National Genomics Data Center in 2020. *Nucleic Acids Res* **48**, D24-d33 (2020).
40. Shi, W. *et al.* gcMeta: a Global Catalogue of Metagenomics platform to support the archiving, standardization and analysis of microbiome data. *Nucleic Acids Res* **47**, D637-d648 (2019).
41. Xiao, S.Z. *et al.* Increased interactivity and improvements to the GigaScience database, GigaDB. *Database (Oxford)* **2019**(2019).
42. Song, S. *et al.* The Global Landscape of SARS-CoV-2 Genomes, Variants, and Haplotypes in 2019nCoV. *Genomics Proteomics Bioinformatics* **18**, 749–759 (2020).
43. World Health Organization. WHO Coronavirus (COVID-19) Dashboard.
44. Africa CDC. Africa Pathogen Genomics Initiative (PGI). Vol. 2021 (2021).
45. Paul, P. *et al.* Genomic Surveillance for SARS-CoV-2 Variants Circulating in the United States, December 2020-May 2021. *MMWR Morb Mortal Wkly Rep* **70**, 846–850 (2021).
46. Brown, K.A. *et al.* S-Gene Target Failure as a Marker of Variant B.1.1.7 Among SARS-CoV-2 Isolates in the Greater Toronto Area, December 2020 to March 2021. *Jama* **325**, 2115–2116 (2021).

Figures

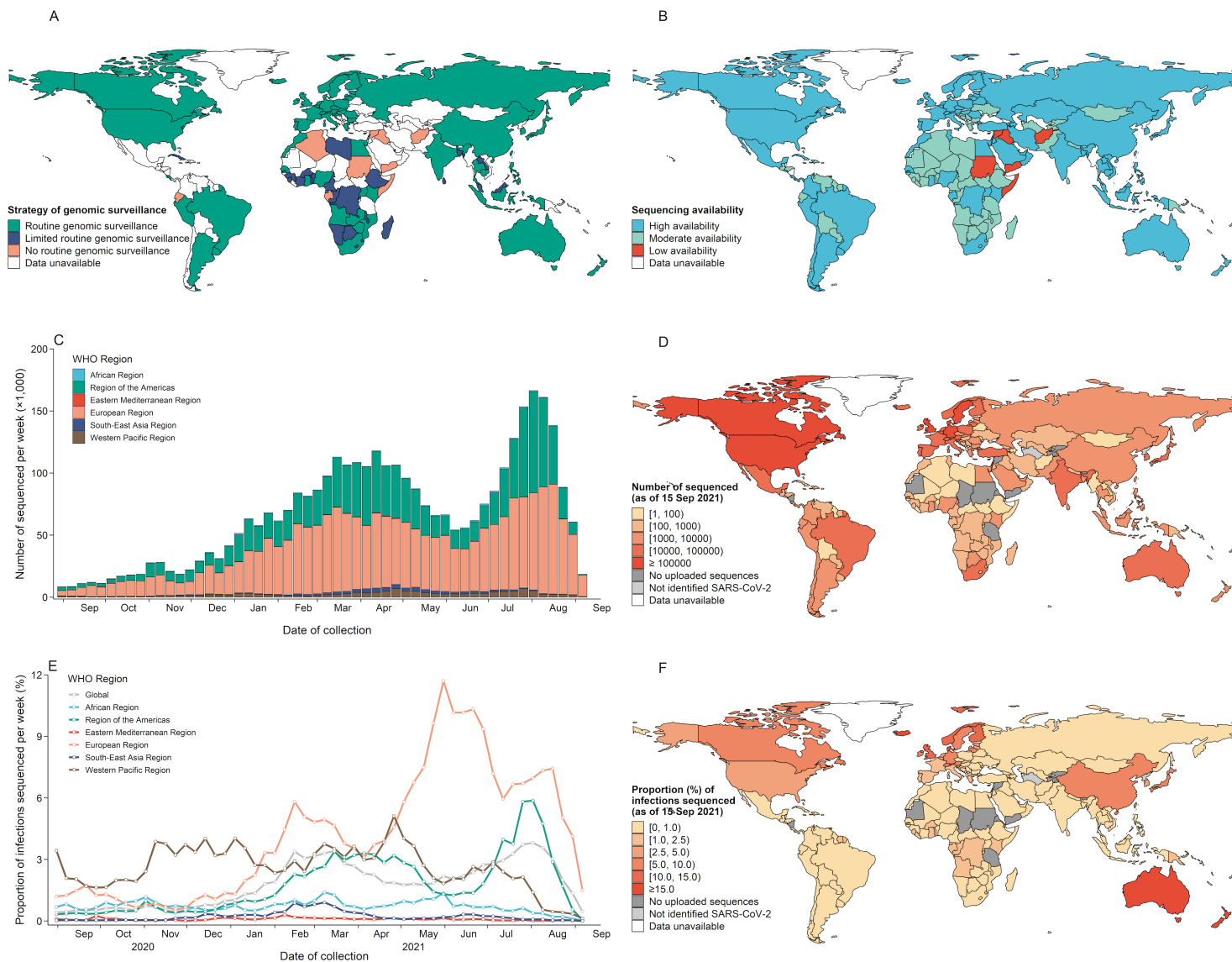


Figure 1

Global SARS-CoV-2 genomic surveillance, sequencing availability, and publicly deposited genomic data. (A) The global distribution of three strategies of SARS-CoV-2 genomic surveillance. (B) The global availability of SARS-CoV-2 sequencing, countries with a high level of availability represent the ability to perform in-country SARS-CoV-2 sequencing alone. (C) The weekly number of publicly deposited SARS-CoV-2 genomic data by region. (D) Cumulative number of publicly deposited SARS-CoV-2 genomic data by countries as of 15 September 2021. (E) The weekly proportion of infections sequenced by region. (F) Cumulative proportion of infections sequenced by countries as of 15 September 2021, defined as the proportion of cumulative isolates sequences to the cumulative confirmed cases. The number of sequences for the most recent weeks might be incomplete due to a time delay between collecting specimens and sequencing submission. Data unavailable, includes those locations that not belonged to 194 Member States or had not applicable data.

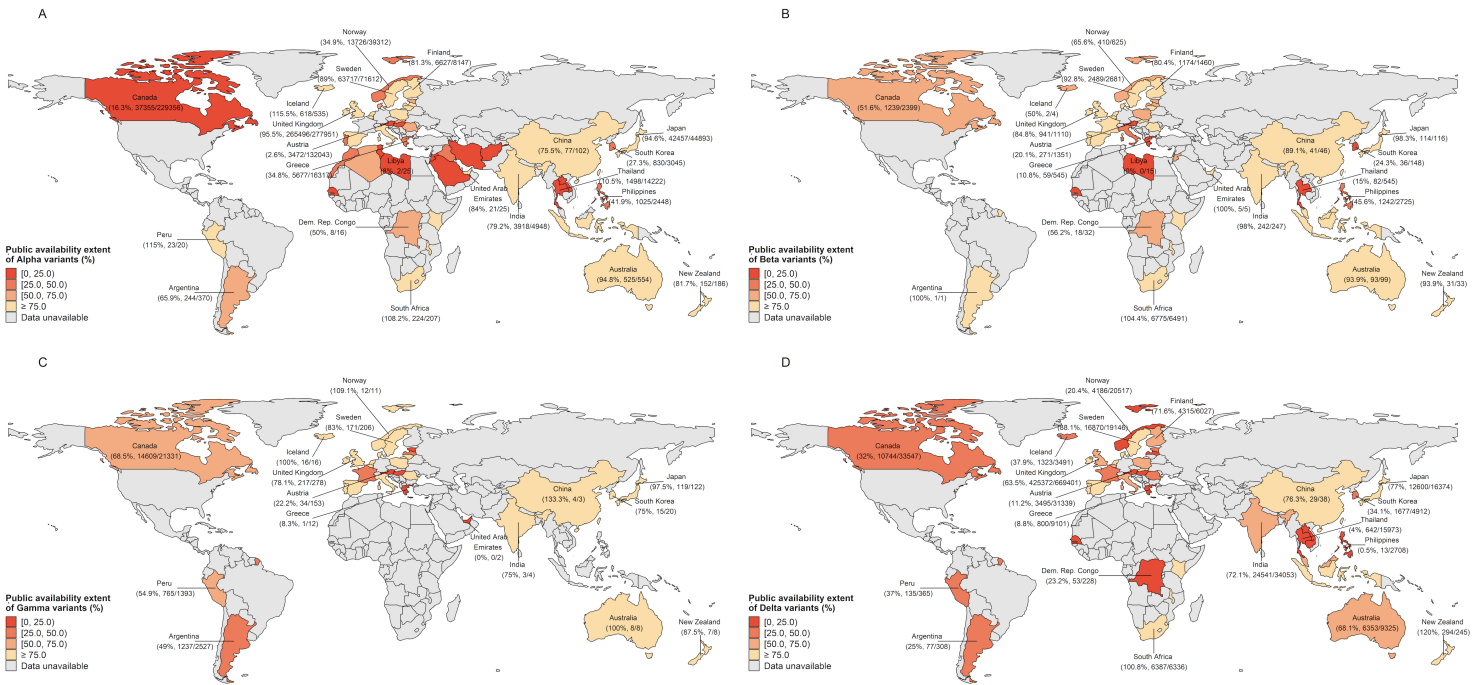


Figure 2

The public availability extent of SARS-CoV-2 genomic data to publicly repositories by country and the variants of concern. The public availability extent was defined as the ratio of the cumulative number of variants in publicly repositories to the official reported number of variants within the same period. The ratio of exceed 100% in some countries might be due to the delay in the official report of the sequencing results or the incomplete official reporting system. In view of the availability of official data, the cumulative number of variants in different countries corresponds to different time periods, with detailed information in Table S6. Sequence without date of specimen collection in publicly repositories is not included in our analysis. The variant data for China only includes cases reported by mainland China. The officially reported number of Alpha variants might contain those confirmed by the PCR screen assay. The values beneath the country names show numbers of cumulative variant as of one specific week: variants in publicly repositories/official reported variants.

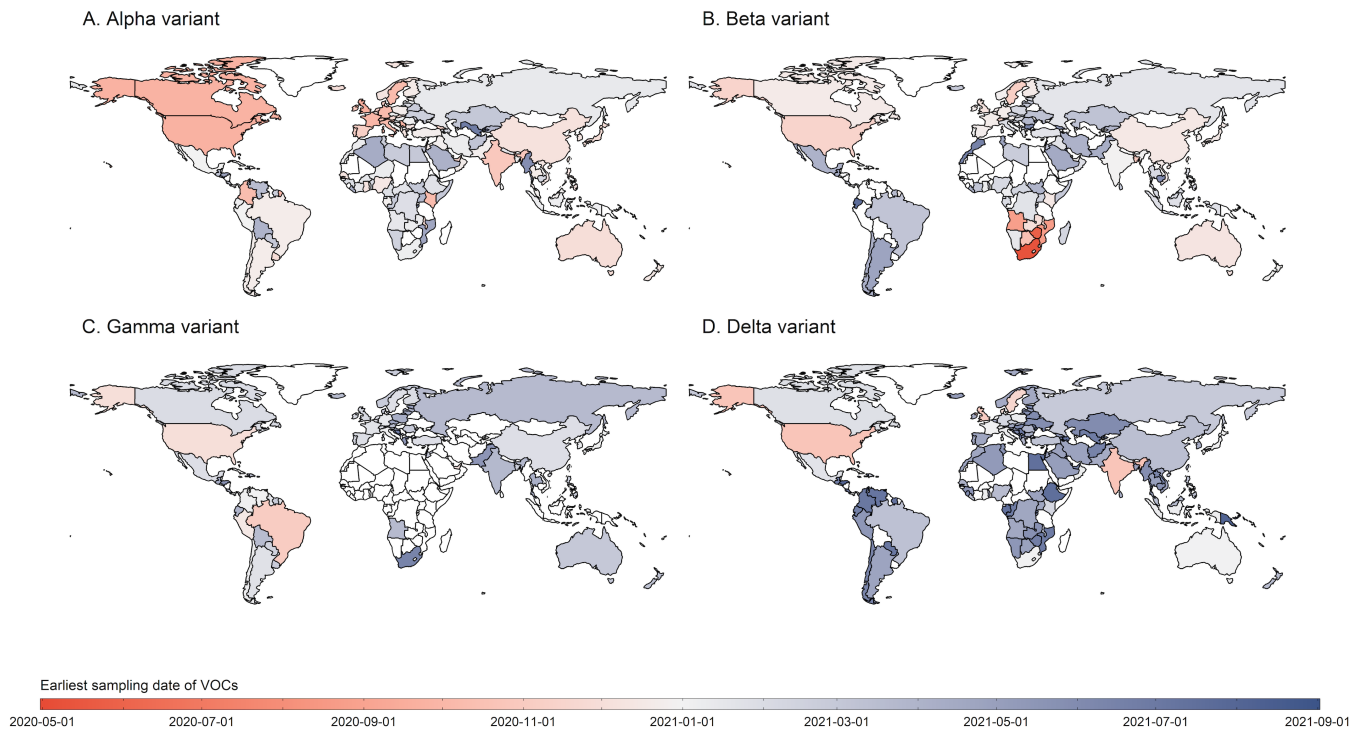


Figure 3

The earliest identification of Alpha, Beta, Gamma, and Delta variant in each country. If information about the earliest sampling date was unavailable but that of the earliest reporting date was available, we extrapolated the sampling date by using a fixed three-week lag from sample collection to reporting. Countries with darker red indicate earlier samples, and with darker blue refers to later samples. The white areas represent countries with variants unreported or data unavailable.

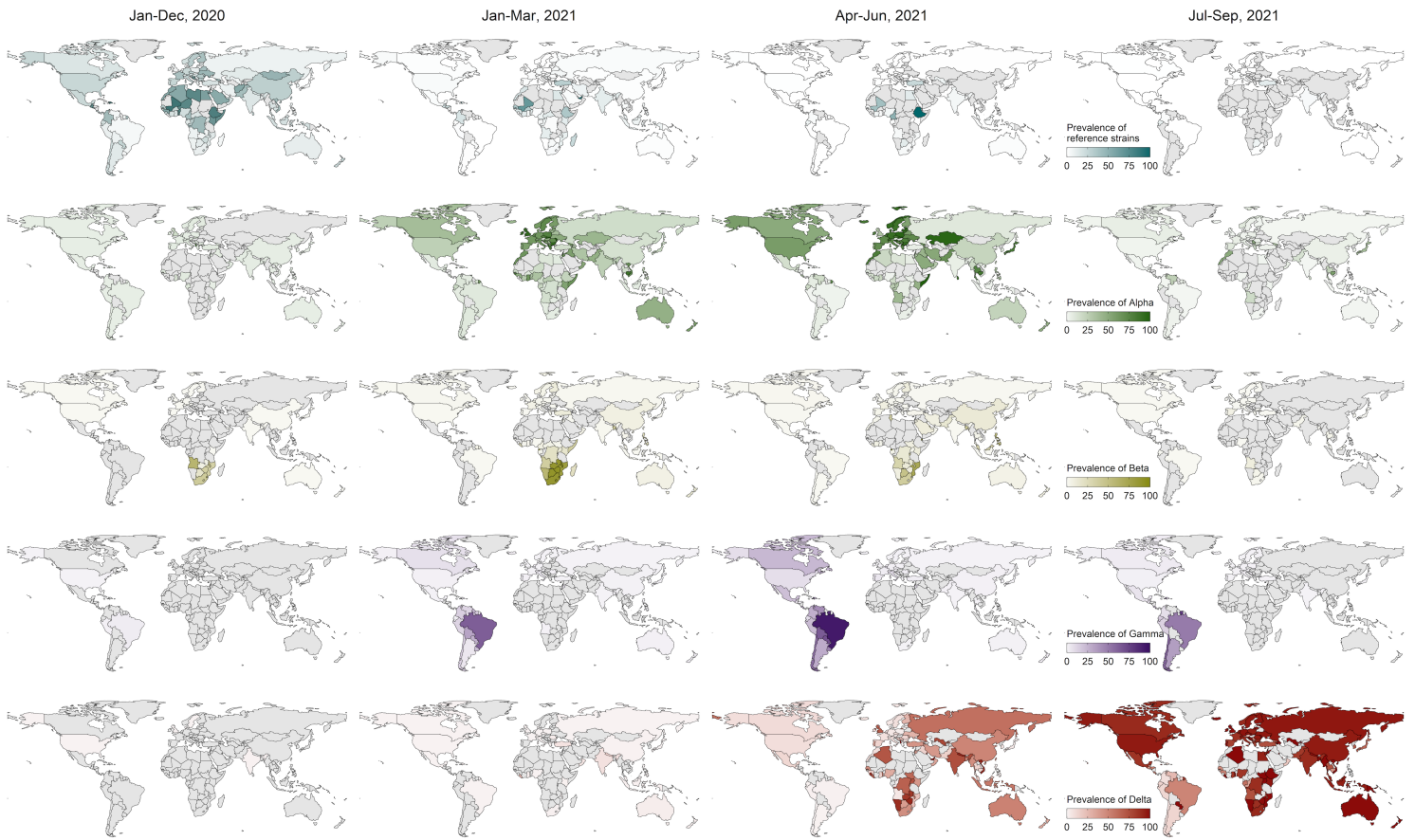


Figure 4

The prevalence and temporal dynamics of reference strains and four SARS-CoV-2 VOCs by country. The date presented in the top refers to the range of date of specimen collection. The prevalence was defined as the proportion of the strain number (reference strains or variants) to the total number of sequences generated in the same unit of time. Reference strains includes lineage A, A.1, B, and B.1; the sub-lineages of four VOCs are aggregated with the parent lineages. The grey areas represent countries with no COVID-19 epidemic, or no performing sequencing, or no uploading genomic data to public repositories.

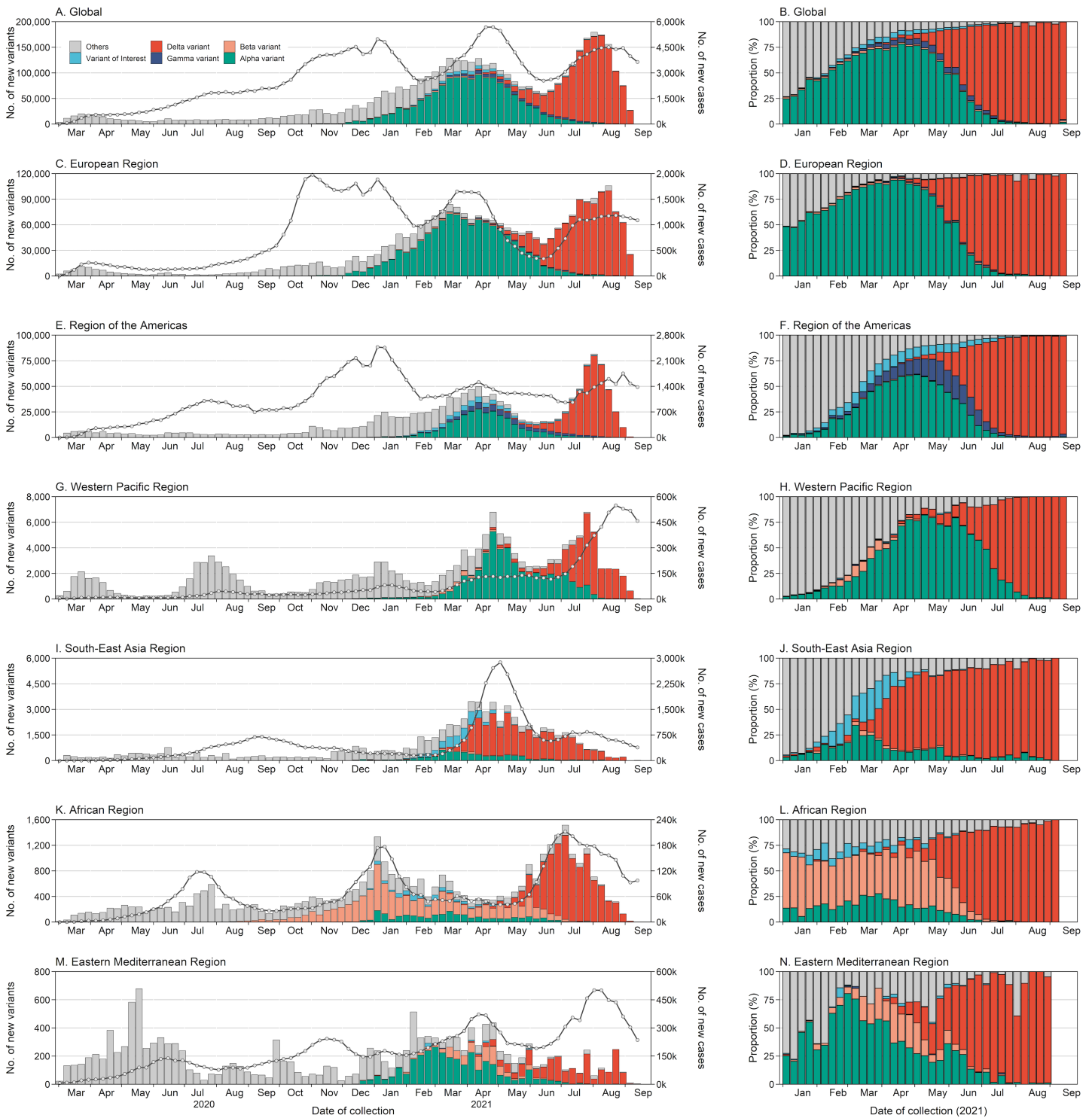


Figure 5

The number and proportion of SARS-CoV-2 variants by region and time. The line and point in the left figure correspond to the y-axis on the right. The sub-lineages of four VOCs are aggregated with the parent lineages; designated Variants of Interest (VOIs) included lineage B.1.525, B.1.526, B.1.617.1, C.37, B.1.621 and their sub-lineages; other lineages included reference strains and other variants. Data used here are derived from the publicly repositories and aggregated dataset, with priority given to the one with highest sequenced numbers in a specific week.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Appendix.docx](#)
- [Acknowledgementtable.pdf](#)