

# Time-Series Trend of Pandemic SARS-CoV-2 Variants Visualized Using Batch-Learning Self-Organizing Map for Oligonucleotide Compositions



RESEARCH PAPER

TAKASHI ABE

RYUKI FURUKAWA

YUKI IWASAKI

TOSHIMICHI IKEMURA

*\*Author affiliations can be found in the back matter of this article*

][ubiquity press

## ABSTRACT

To confront the global threat of coronavirus disease 2019, a massive number of the severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) genome sequences have been decoded, with the results promptly released through the GISAID database. Based on variant types, eight clades have already been defined in GISAID, but the diversity can be far greater. Owing to the explosive increase in available sequences, it is important to develop new technologies that can easily grasp the whole picture of the big-sequence data and support efficient knowledge discovery. An ability to efficiently clarify the detailed time-series changes in genome-wide mutation patterns will enable us to promptly identify and characterize dangerous variants that rapidly increase their population frequency. Here, we collectively analyzed over 150,000 SARS-CoV-2 genomes to understand their overall features and time-dependent changes using a batch-learning self-organizing map (BLSOM) for oligonucleotide composition, which is an unsupervised machine learning method. BLSOM can separate clades defined by GISAID with high precision, and each clade is subdivided into clusters, which shows a differential increase/decrease pattern based on geographic region and time. This allowed us to identify prevalent strains in each region and to show the commonality and diversity of the prevalent strains. Comprehensive characterization of the oligonucleotide composition of SARS-CoV-2 and elucidation of time-series trends of the population frequency of variants can clarify the viral adaptation processes after invasion into the human population and the time-dependent trend of prevalent epidemic strains across various regions, such as continents.

## CORRESPONDING AUTHORS:

### Takashi Abe

Smart Information Systems,  
Faculty of Engineering, Niigata  
University, Niigata-ken 950-  
2181, Japan

[takaabe@ie.niigata-u.ac.jp](mailto:takaabe@ie.niigata-u.ac.jp)

### Toshimichi Ikemura

Department of Bioscience,  
Nagahama Institute of  
Bio-Science and Technology.  
Shiga-ken 526-0829, Japan

[t\\_ikemura@nagahama-i-bio.ac.jp](mailto:t_ikemura@nagahama-i-bio.ac.jp)

## KEYWORDS:

COVID-19; SARS-CoV-2;  
Oligonucleotide composition;  
Batch-Learning Self-Organizing  
Map (BLSOM); Unsupervised  
explainable machine learning;  
Time-series trend

## TO CITE THIS ARTICLE:

Abe, T, Furukawa, R, Iwasaki,  
Y and Ikemura, T. 2021.  
Time-Series Trend of Pandemic  
SARS-CoV-2 Variants Visualized  
Using Batch-Learning  
Self-Organizing Map for  
Oligonucleotide Compositions.  
*Data Science Journal*, 20: 29,  
pp. 1–12. DOI: <https://doi.org/10.5334/dsj-2021-029>

## INTRODUCTION

The severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) has spread rampantly worldwide since it was first reported in December 2019, and the momentum of its spread is still increasing (WHO. 2020). To address the SARS-CoV-2 pandemic in detail, genome sequencing has been performed on a global scale and published by GISAID (Elbe et al. 2017), the SARS-CoV-2 genome database, having more than 780,000 viral sequences as of March 2021 (<https://www.gisaid.org/>). SARS-CoV-2 is an RNA virus with a fast evolutionary rate that has already been classified into eight clades by GISAID, and epidemics caused by new variant have been known to occur (Benvenuto et al. 2020; Gorbalenya et al. 2020; Sun et al. 2020; Hu et al. 2021; Kirby 2021; Wang et al. 2021). Because the number of registered genome sequences is increasing explosively, it has become difficult to cope with the current and future situation using only the conventional phylogenetic tree method based on multiple sequence alignment, which requires an enormous amount of computation time for a massive number of sequences. Therefore, it is imperative to develop a sequence alignment-free method that will enable us to easily grasp the whole picture of the big-sequence data and support efficient knowledge discovery from them.

By focusing on the frequency of short oligonucleotides (e.g., tetra- and penta-nucleotides) in a large number of genomic fragments (e.g., 10 kb) derived from a wide variety of species, we have developed an unsupervised explainable AI (batch-learning self-organizing map; BLSOM), which enables separation (self-organization) of the genomic sequences by species and phylogeny and explains the causes that contribute to this separation (Abe et al. 2003 ). In the analysis of genomic fragments of a wide range of microbial genomes, over 5 million sequences can be separated by phylogenetic groups with high accuracy (Abe et al. 2020).

In a prior analysis of all influenza A strains, viral genomes were separated (self-organized) by host animals based only on the similarity of the oligonucleotide composition, although no host information was provided during BLSOM learning (Iwasaki et al. 2011). On a single map, all viral sequences could be separated, and notably, BLSOM is an explainable AI that can explain diagnostic oligonucleotides, which contribute to host-dependent clustering. When studying the 2009 swine-derived flu pandemic (H1N1/2009), we could detect directional time-series changes in oligonucleotide composition because of possible adaptations to the new host, namely humans (Iwasaki et al. 2011), showing that near-future prediction was possible, albeit partially (Iwasaki et al. 2013).

We have previously revealed lineage-specific oligonucleotide compositions for a wide range of virus lineages and established a method to identify and classify viral-derived sequences in tick intestinal metagenomic sequences (Qiu et al. 2019). In the case of SARS-CoV-2, we analyzed time-series changes in mono- and oligo-nucleotide compositions and found their time-dependent directional changes that are thought to be adaptive for growth in humans, which allowed us to predict candidates of advantageous mutations for growth in human cells (Ikemura et al. 2020; Wada, Wada & Ikemura. 2020; Iwasaki Abe & Ikemura. 2021). Furthermore, we recently performed BLSOM analysis on di- to penta-nucleotide compositions in approximately 150,000 SARS-CoV-2 genomes. Because the accuracy of separation by clade increased as the oligonucleotide length increased, in this report, we present the BLSOM results for the pentanucleotide composition. BLSOM could serve as a powerful tool for comprehensive characterization of the oligonucleotide composition of SARS-CoV-2 and time-series trends of prevalent epidemic strains across various regions, such as continents.

## METHODS

### SARS-COV-2 GENOME SEQUENCES

The full-length genome sequences of SARS-CoV-2 were downloaded from the GISAID database on November 4, 2020. The total number of sequences was 170,190. The full length of the SARS-CoV-2 genome reference sequence (strain name: Wuhan-Hu-1, accession number: MN908947.3), which includes 5' and 3' untranslated regions (UTRs) and polyA tail, is 29.9 kb. To analyze more genome data, after removing the poly (A)-tail sequences, we set the minimum threshold length to 27 kb, which includes a major part of coding sequence.

Pentanucleotide frequencies and odds ratios were used in the present study. The pentanucleotide odd ratios (observed/expected values) were calculated using the formula  $P_{VWXYZ} = f_{VWXYZ} / f_V f_W f_X f_Y f_Z$ , where  $f_V$ ,  $f_W$ ,  $f_X$ ,  $f_Y$  and  $f_Z$  denote the frequencies of mononucleotides V, W, X, Y and Z, respectively, and  $f_{VWXYZ}$  denotes the frequency of pentanucleotide VWXYZ (Karlin et al. 1998).

## BLSOM

Kohonen's self-organizing map (SOM), an unsupervised neural network algorithm, is a powerful tool for clustering and visualizing high-dimensional complex data on a two-dimensional map (Kohonen, 1990; Kohonen et al., 1996). We modified the conventional SOM for genome informatics on the basis of batch learning, aiming to make the learning process and the resulting map independent of the order of data input (Kanaya et al. 2001; Abe et al. 2003). The newly developed SOM, BLSOM, is suitable for high-performance parallel computing and, therefore, for big data analysis. The initial weight vectors were defined using principal component analysis (PCA), based on the variance-covariance matrix, rather than by using random values. The weight vectors ( $w_{ij}$ ) were arranged in a two-dimensional lattice denoted by  $i$  ( $= 0, 1, \dots, I-1$ ) and  $j$  ( $= 0, 1, \dots, J-1$ ) and were set and updated as described previously (Kanaya et al. 2001; Abe et al. 2003). A BLSOM program suitable for PC cluster systems is available on our website (<http://bioinfo.ie.niigata-u.ac.jp/?BLSOM>). After constructing BLSOM and its 3-D view explained in the text, we first assigned the lattice point that has high number of sequences in each continent to the representative point of the continent and manually defined the lattice points surrounding the representative point as subclusters.

## RESULTS AND DISCUSSION

### BLSOM FOR PENTANUCLEOTIDE COMPOSITION AND THEIR ODDS RATIO

It should be mentioned here that SARS-CoV-2 genomes have changed their mononucleotide composition during the course of the epidemic in humans, reducing C and increasing U, regardless of clade (Mercatelli et al. 2020; Wada, Wada & Ikemura. 2020; Iwasaki Abe & Ikemura 2021), a process which is thought to be caused by the APOBEC family enzymes (Mangeat et al. 2003; Simmonds 2020). It should also be noted here that GISAID clade was defined by a nomenclature system developed by the GISAID group and divided into seven clades, including S, L, V, G, GH, GR and GV, based on marker mutations by November 2020. The clade division was initially S and L during the early epidemic stage, but L was further divided into V and G, and then later, G was divided into GH, GR and GV. The marker mutations of these clades include NS8-L84S for clade S, NSP6-L37F and NS3-G251V for clade V, and S-D614G for clade G. In addition to clade G, NS3-Q57H, N-G204R and S-A222V mutations define the clades GH, GR and GV, respectively (Elbe et al. 2017, <https://www.gisaid.org/>). Considering the clade-independent tendency primarily caused by the APOBEC enzymes (Simmonds 2020), we performed BLSOM analysis of not only the pentanucleotide composition but also their odds ratio, which can reduce the effects caused by changes in mononucleotide composition. Additionally, to check the robustness of sequence accuracy, we used datasets with different sequence accuracies: 167,905 sequences with less than 10% unknown nucleotides other than ATGCs in the genome sequence and 130,753 sequences with less than 1% unknown nucleotides; for each sequence dataset, the number of cases by region and clade is shown in [Table 1](#).

First, we constructed BLSOM for sequences with less than 10% unknown nucleotides, using the pentanucleotide composition and their odds ratios ([Figure 1A](#) and [B](#)). BLSOM utilizes unsupervised machine learning, and the genome sequences are clustered (self-organized) on a two-dimensional plane, based only on the difference in the vector data in a  $1024 (=4^5)$ -dimensional space. Lattice points that include sequences from more than one clade are indicated in black, those that contain no genomic sequences are indicated by blank, and those containing sequences from a single clade are indicated in the color representing the clade. The odds ratio ([Figure 1B](#)) gave more accurate separations (a smaller percentage of black grid points), possibly by excluding effects owing to the clade-independent time-series change in the mononucleotide composition (Iwasaki Abe & Ikemura. 2021), which affected all SARS-CoV-2 clades. Even for the sequences with low-sequence accuracy, clade-dependent separation occurs, allowing us to understand characteristics of the oligonucleotide composition that are

<b>(A) NUMBER OF SEQUENCES WITH LESS THAN 10% UNKNOWN NUCLEOTIDES</b>								
<b>CLADE\ CONTINENT</b>	<b>ASIA</b>	<b>EUROPE</b>	<b>NORTH AMERICA</b>	<b>OCEANIA</b>	<b>AFRICA</b>	<b>SOUTH AMERICA</b>	<b>UNKNOWN</b>	<b>TOTAL</b>
S	794	1,860	3,449	664	110	74	0	6,951
L	823	3,196	600	65	4	11	0	4,699
V	247	4,687	402	253	13	23	0	5,625
G	979	20,928	6,568	1,106	1,141	461	0	31,183
GH	2,058	10,325	23,916	964	232	176	0	37,671
GR	2,657	42,888	5,251	11,135	1,632	1,129	0	64,692
GV	3	12,229	3	14	0	0	0	12,249
O	2,220	1,127	553	531	60	25	0	4,516
Non-human host	35	247	19	0	1	4	13	319
#Total	9,816	97,487	40,761	14,732	3,193	1,903	13	167,905

<b>(B) NUMBER OF SEQUENCES WITH LESS THAN 1% UNKNOWN NUCLEOTIDES</b>								
<b>CLADE\ CONTINENT</b>	<b>ASIA</b>	<b>EUROPE</b>	<b>NORTH AMERICA</b>	<b>OCEANIA</b>	<b>AFRICA</b>	<b>SOUTH AMERICA</b>	<b>UNKNOWN</b>	<b>TOTAL</b>
S	731	1,047	3,056	466	71	58	0	5,429
L	760	1,964	549	49	2	10	0	3,334
V	228	3,036	366	207	10	17	0	3,864
G	877	15,200	5,071	858	634	300	0	22,940
GH	1,923	8,365	19,014	717	191	150	0	30,360
GR	2,425	32,518	4,549	9,166	1,180	871	0	50,709
GV	3	10,712	3	11	0	0	0	10,729
O	1,824	522	349	415	30	9	0	3,149
Non-human host	30	176	19	0	1	0	13	239
#Total	8,801	73,540	32,976	11,889	2,119	1,415	13	130,753

**Table 1** Number of SARS-CoV-2 genome sequences with less than 10% (A) and less than 1% (B) unknown nucleotides used in this study. Unknown: genome sequences for which continent was not registered.

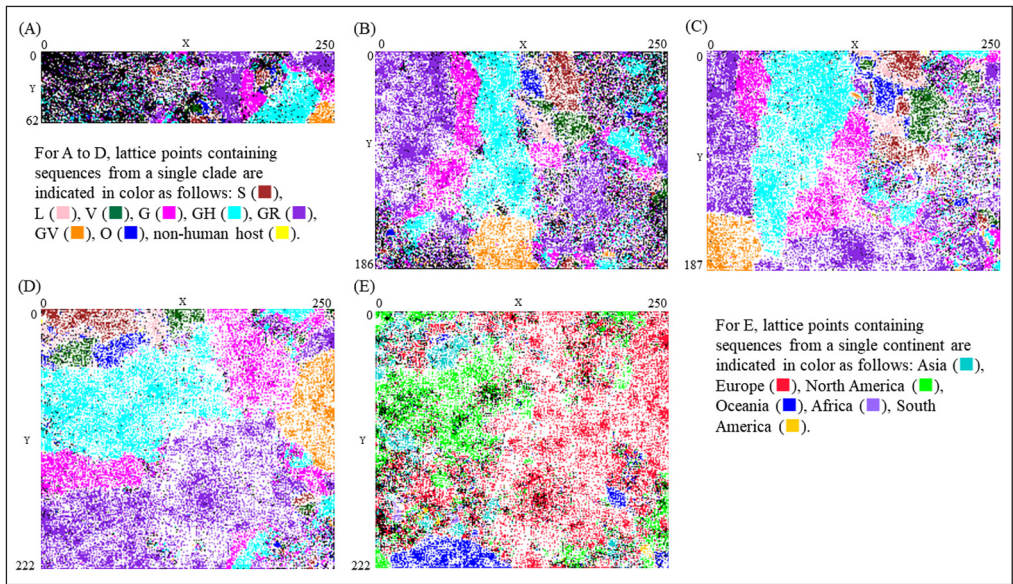
specific to each clade; thus, oligonucleotide-BLSOM is thought to be a robust method. However, it is clear that BLSOMs for sequences with less than 1% unknown nucleotides (Figure 1C and D) gave more accurate separation than those listed in Figure 1A and B, and the highest resolution was obtained for the BLSOM for the odds ratio (Figure 1D).

Clades have been defined by the statistical distribution of phylogenetic distances in tree construction based on multiple sequence alignments (Han et al. 2019; Tang et al. 2020), whereas BLSOM is a sequence alignment-free analysis that is suitable for the analysis of massive data. Because sequences at different locations on BLSOM have different oligonucleotide compositions, clustering according to clades means that sequences belonging to different clades have different oligonucleotide combinations, that is, differential combinations of mutations.

### 3D DISPLAY OF THE DATA FOR DIFFERENT CONTINENTS

Using BLSOM (Figure 1D) for the pentanucleotide odds ratio, Figure 1E examines the classification according to four continents (Asia, Europe, North America, and Oceania) that have very large numbers of sequences and thus selected as the main epidemic continents. Here, the lattice points containing sequences of different continents are displayed in black, and those containing only sequences of a single continent are displayed in the color specifying each continent. Although not as clear as clade-dependent separations, regional differences have been observed, which should reflect differential shares of prevalent variants among continents. However, it is apparently difficult to obtain sufficient information from the results shown in





**Figure 1** BLSOM for pentanucleotide usage. (A) Pentanucleotide composition and (B) their odds ratio for sequences with less than 10% unknown nucleotides. (C) Pentanucleotide composition and (D) their odds ratio for sequences with less than 1% unknown nucleotides. Lattice points that include sequences from more than one clade are indicated in black, those that contain no genomic sequences are indicated by blank, and those containing sequences from a single clade are indicated in color as follows: S (■), L (■), V (■), G (■), GH (■), GR (■), GV (■), O (■), non-human host (■). (E) Distribution of sequences by continent on the BLSOM with the pentanucleotide odds ratio. Lattice points that include sequences from more than one continent are indicated in black, those that contain no genomic sequences are indicated by blank, and those containing sequences from a single continent are indicated in color as follows: Asia (■), Europe (■), North America (■), Oceania (■), Africa (■), South America (■).

Figure 1E alone. BLSOM is equipped with various visualization tools for analysis results; therefore, we next showed the number of sequences belonging to each lattice point with a 3D display.

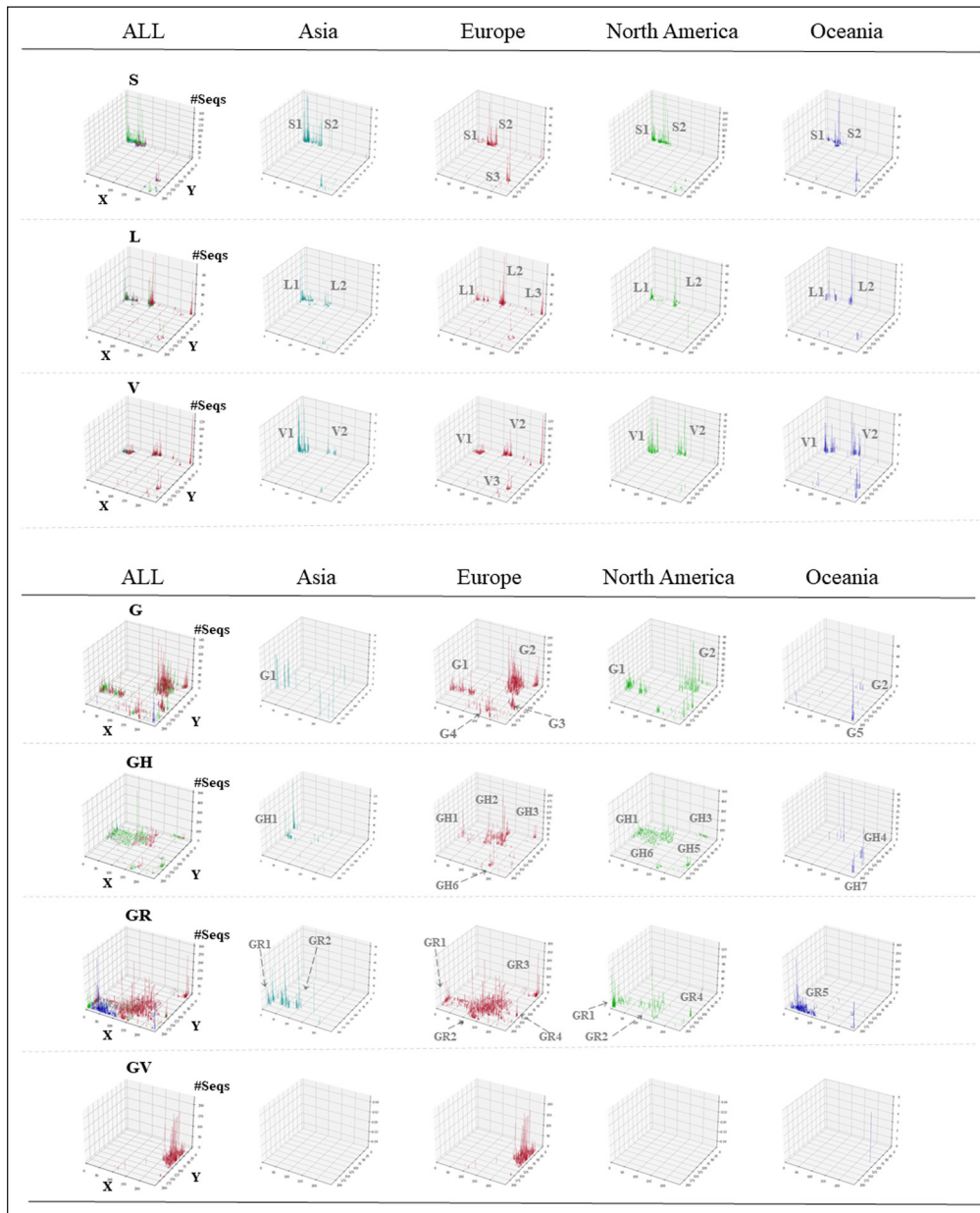
Again, using the BLSOM shown in Figure 1D, Figure 2 shows the number of sequences belonging to each lattice point for each clade in each continent as a vertical bar, which is colored by continent, as shown in Figure 1E. Looking laterally at a particular clade, each clade consists of several subclusters, each consisting of several high peaks surrounded by many low peaks. Different subclusters observed in each clade are distinguished by numbering in each figure, but if they are located in the same zone on BLSOM, the same number is given even if they are of different continents. Looking vertically at a particular continent, sequences of different subclusters of different clades exist in different amounts, and some subclusters are only in a particular continent, that is, the prevalent variants for each continent can be visualized in an easy-to-understand manner. In Supplementary Figure S1, the data shown in Figure 2 are displayed in 2D, and referring to the quantitative results in Figure 2, we defined sequences attributed to each subcluster in each clade.

**TIME-SERIES ANALYSIS**

The fact that sequences belonging to one clade were clearly separated on BLSOM indicates the importance of subdivision of each clade, and the separation on BLSOM is thought to be a good indicator of this subdivision. To further examine the biological significance of the subclusters of each clade on BLSOM, we visualized the number of sequences collected in each month in each region as a vertical bar differentially colored according to clade (Figure 3). Looking laterally at a continent, the time-series quantitative changes among different clades or different subclusters of one clade are clear. Looking at the results for a particular collection month for different continents longitudinally, quantitative changes among different clades or different subclusters of one clade are again clear, depending on the continent.

Next, for each clade in each continent, we quantitatively analyzed the time-series changes in the proportion of its subclusters using a 100% stack bar graph (Figure 4). The percentages of sequences in different subclusters are distinguished by different colors, and when the total number of sequences for a certain month is more than 100, the data for that month are indicated by a thick horizontal bar. We focused mainly on such months.

In the clade S/L/V detected in the early stage of the epidemic (December 2019– March 2020), three major subclusters of each clade were observed and distinguished by suffix numbers, and most sequences belonged to the two subclusters: S1/L1/V1 and S2/L2/V2. In Asia, many sequences belonging to S1/L1/V1 were detected in December 2019, but in Europe and other regions, S2/L2/V2 were more abundantly detected in March and April 2020 than S1/L1/V1, and the proportion became more pronounced in April than in March. In March and April in Europe, a remarkable number of sequences belonging to S3/L3/V3 were also detected, showing three different variants prevalent at the beginning of the epidemic in Europe. Far fewer than 100



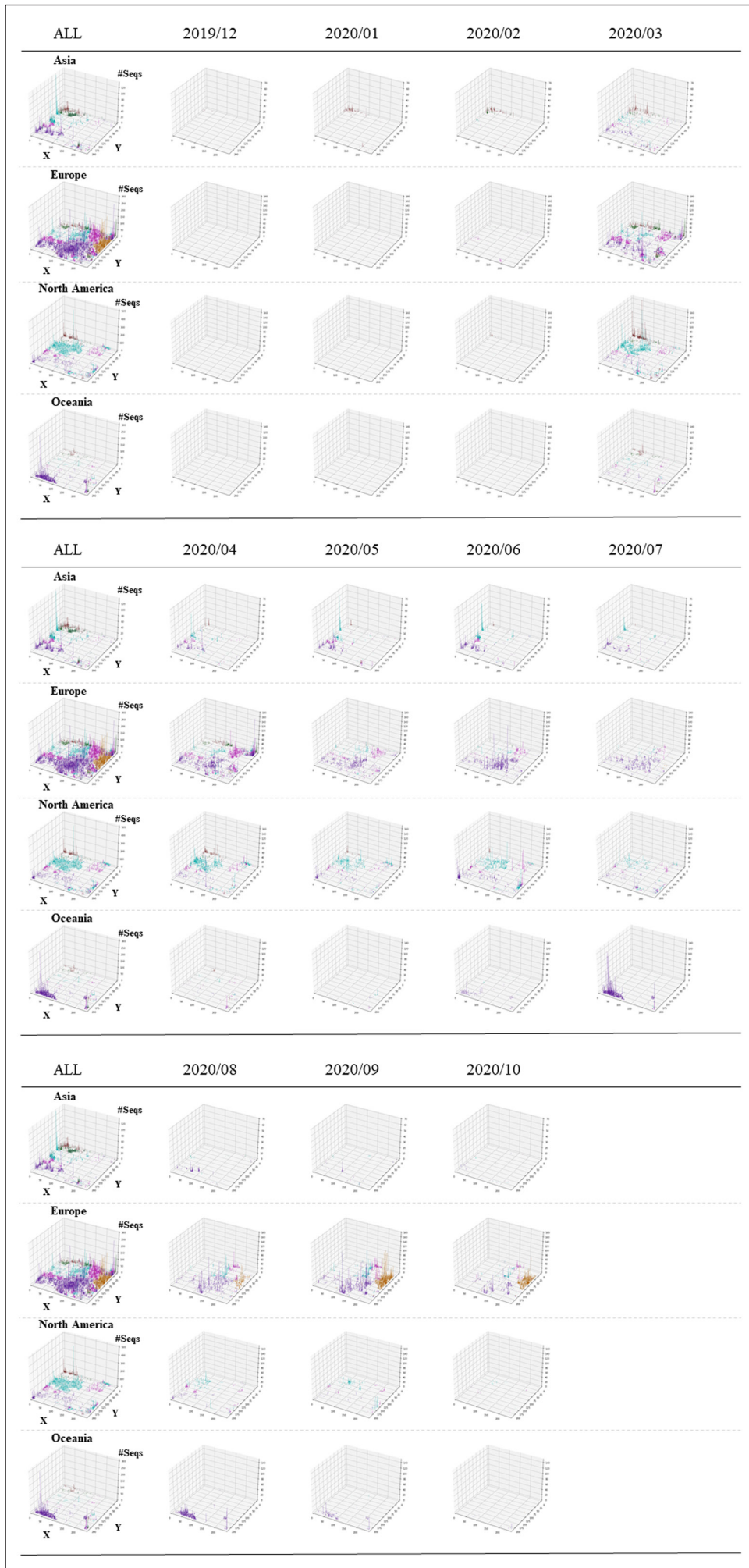
**Figure 2** 3D display of viral classification by clade and continent. The Z-axis corresponds to the number of sequences attributed to each lattice point. Results for all continents are shown in the ALL panel for each clade. In clades G, GH, GR and GV, lattice points where less than 5 sequences exist are not shown. The vertical bars for individual continents are distinguished by the following colors: Asia (■), Europe (■), North America (■), Oceania (■). Different subclusters are given suffix numbers.

sequences were detected after May; sequences belonging to S1/L1/V1 were mainly detected in Asia and those belonging to S2/L2/V2 were shown in other regions, presenting differential trends in prevalent variants among continents.

For clade G, which started the epidemic in Europe in February, we defined five subclusters. In February, roughly equal amounts of sequences belonging to G1 and G2 were detected in Europe and North America, but as the epidemic progressed, those belonging to G2 were mainly detected in Europe, whereas those belonging to both G1 and G2 were prevalent in North America. In Asia, only sequences belonging to G1 were detected; in Oceania, those belonging to G2 accounted for about 10% in the early stage, but afterward, those belonging to Oceania-specific G5 accounted for the majority.

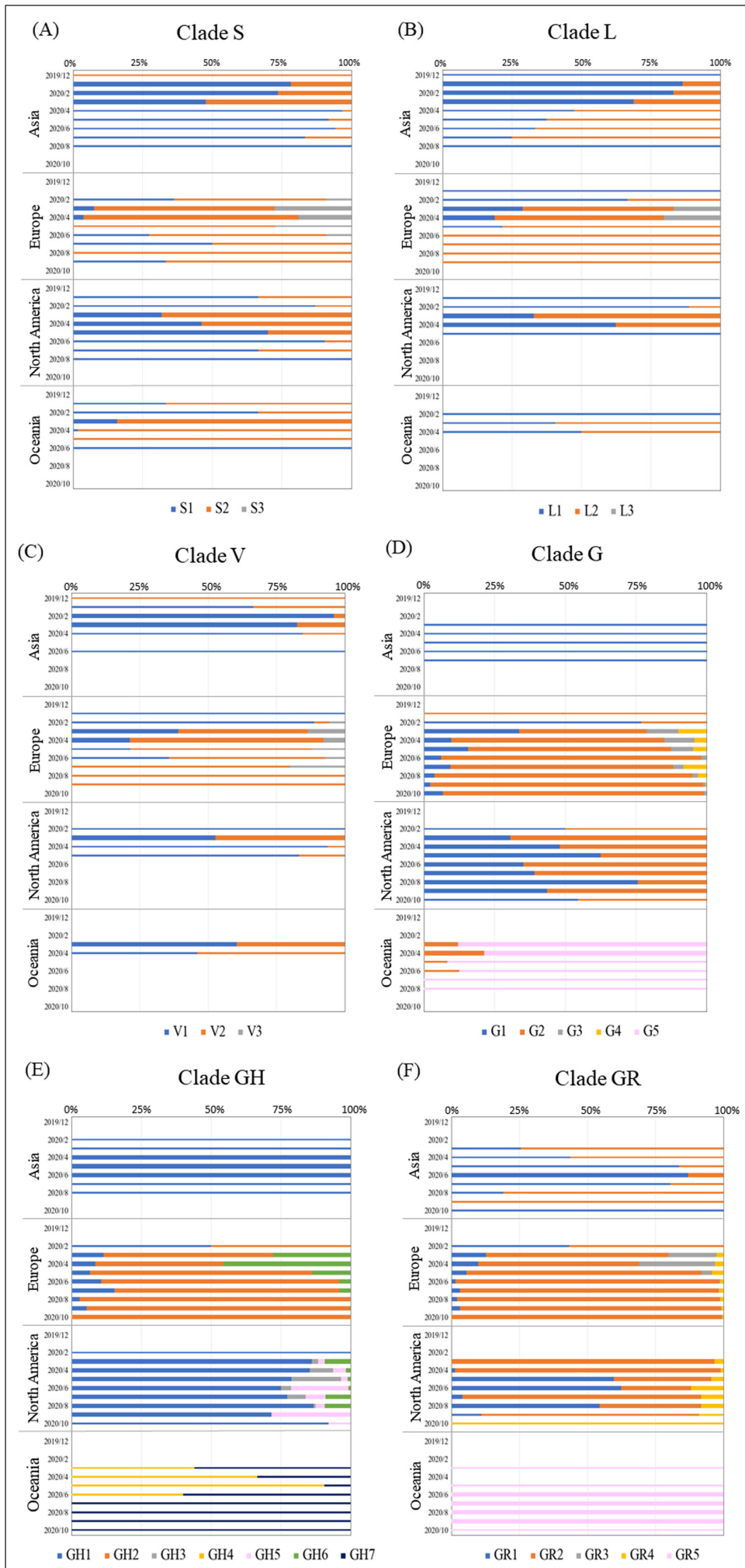
For GH, we defined seven subclusters, including GH1 and GH2, which dominated in North America and Europe, respectively. In North America, in addition to GH1, several months contain approximately 20% of the sequences belonging to GH3, GH5, and GH6. In Asia, only GH1 has been detected. In Oceania, only GH4 and GH7, which were specific to this region, were detected; initially, GH4 was dominant, but after July, GH7 was primarily detected.

For GR, we defined five subclusters, including GR1 and GR2, which dominated in North America and Europe, respectively. Moreover, in Europe, GR1 was detected to the same extent as GR2 in February, but as the epidemic progressed, GR2 began to predominate. In North America, the occupancy of GR1 and GR2 varied to some extent depending on the collection month. In Asia, GR1 was mainly detected, and in Oceania, only region-specific subclusters have been detected.



**Figure 3** 3D display of temporospatial changes. The Z-axis corresponds to the number of sequences attributed to each lattice point. Results for all collection months are shown in the ALL panel for each continent. The vertical bars for individual clades are distinguished by the following colors: S (red), L (pink), V (green), G (purple), GH (cyan), GR (blue), GV (orange).





**Figure 4** Analysis of 100% stack bar graph for time-series transition in each continent for each subcluster in clades S (A), L (B), V (C), G (D), GH (E), and GR (F). The colors of each subcluster are indicated at the bottom of each figure. The results for months with more than 100 sequences are shown as thick horizontal bars. The number of sequences used in this analysis is given in Supplementary Table S1.



These temporospatial changes in subclusters show that the subcluster is the separation (self-organization) that reflects biological significance and is fundamental information for understanding the overall picture of the SARS-CoV-2 variants.

## BIOLOGICAL MEANINGS OF BLSOM SEPARATION

Change in oligonucleotide composition is strongly influenced by changes in mononucleotide composition; the CU mutation in SARS-CoV-2 caused by APOBEC is well known (Simmonds, 2020). However, in a time-series study, we have found many changes that cannot be explained by mononucleotide changes (Iwasaki et al. 2021). If a mutation occurs that alters protein function and clearly increases infectivity or growth rate, the mutation will rapidly increase its frequency in the viral population and lead to the formation of a new clade, resulting in BLSOM separation. Notably, there are many synonymous mutations that have rapidly increased their frequencies, and detailed time-series analyses of their population frequency showed that some synonymous mutations appear not to be neutral (Wada et al. 2020). Concerning oligonucleotides such as pentanucleotides, some are expected to bind to host proteins or RNAs, and oligonucleotides adapting well to the host factors in human cells may differ from those adapting to natural hosts (e.g., bat). When functionally advantageous mutations including synonymous ones occur, they are thought to lead to the emergence of a new clade and BLSOM separation. At this time, we do not have a clear answer to the actual molecular mechanisms of the possible advantageous mutations, but we are analyzing them as separate studies (Wada et al. 2020; Ikemura et al. 2021). In a previous study, we assigned mutations that contribute to the separation on BLSOM; these diagnostic mutations including synonymous ones are found not only in the spike protein gene but also in many other genes (Ikemura et al. 2020 & 2021).

## CONCLUSION AND PERSPECTIVES

Based on the phylogenetic tree construction by multiple sequence alignments, GISAID has defined seven clades of SARS-CoV-2, giving a total of eight if clade O corresponding to others is included. However, these classifications are clearly inadequate to understand the current status of SARS-CoV-2 because this RNA virus evolves at a high speed. Using only the oligonucleotide composition of many genomic sequences, the unsupervised machine learning, BLSOM, could separate viral sequences according to not only clades but also subclusters within each clade. The separation (self-organization) that AI can accomplish without any hypothesis or model is thought to be a classification from a new perspective. BLSOM is equipped with various tools that allow us to visualize the analysis results in an easily understandable way and to visualize differences in the number of subcluster sequences among continents ([Figure 2](#)) and their time-series changes ([Figure 3](#)), i.e., the distinct variations in the resulting subclusters depending on the region and the collection time.

Herein, we focused on pentanucleotide composition, but similar separations were obtained for other lengths of oligonucleotides (Ikemura et al. 2020). BLSOM is an explanatory AI that can clarify combinatorial patterns of oligonucleotides that contribute to the separation according to clades and their subclusters. BLSOM is a powerful method for comprehensive characterization of the oligonucleotide composition in a massive number of SARS-CoV-2 genome sequences. Next, it will be important to know the relationship between the strains isolated in clades and their subclusters and the causative mutations. When it comes to oligonucleotides as long as 15-mers, most are only present in one copy in the viral genome; therefore, changes in 15-mer sequences can be directly linked to mutations, and we have already started analysis from this perspective (Ikemura et al. 2020). The implementation of time-series oligonucleotide analysis of variants with rapidly expanding intra-population frequencies has enabled the identification of candidates for advantageous mutations for viral infection and growth in human cells (Wada, Wada & Ikemura 2020).

Phylogenetic methods based on sequence alignment have been widely used in evolutionary studies (Hadfield et al. 2018; Kumar et al. 2018), and these methods are undoubtedly essential for studying the phylogenetic relationships between different viral species and variations in the same virus at the single-nucleotide level. In contrast, AI can analyze a massive number of SARS-CoV-2 sequences at once without difficulty, potentially reaching a level of one million

## ADDITIONAL FILES

The additional files for this article can be found as follows:

- **Supplementary Figure S1.** 2D display of the classification by clade and continent shown in Figure 2. Each subcluster territory is circled by a dotted line. In clades G, GH, GR, and GV, lattice points where less than 5 sequences exist are not shown. The sequences belonging to each territory defined here are used for the analysis in Figure 4. DOI: <https://doi.org/10.5334/dsj-2021-029.s1>
- **Supplementary Table S1.** Sequence number of subdivided clusters in clade for each month by continent. DOI: <https://doi.org/10.5334/dsj-2021-029.s2>

## ACKNOWLEDGEMENTS

We gratefully acknowledge the authors submitting their sequences of the GISAID Database.

## FUNDING INFORMATION

This research was supported by AMED Grant Number JP20he0622033h0001, and by JST, CREST Grant Number JPMJCR20H1, and by KAKENHI Grant Numbers 18K07151 and 20H03140.

## COMPETING INTERESTS

The authors have no competing interests to declare.

## AUTHOR CONTRIBUTIONS

TA and TI conceived and designed the research; TA, RF, and TI performed the research; TA, RF, and YI analyzed the data; and all authors wrote the paper.

## AUTHOR AFFILIATIONS

**Takashi Abe**  [orcid.org/0000-0002-8692-4241](https://orcid.org/0000-0002-8692-4241)

Smart Information Systems, Faculty of Engineering, Niigata University, Niigata-ken 950-2181, Japan

**Ryuki Furukawa**

Smart Information Systems, Faculty of Engineering, Niigata University, Niigata-ken 950-2181, Japan

**Yuki Iwasaki**

Department of Bioscience, Nagahama Institute of Bio-Science and Technology. Shiga-ken 526-0829, Japan

**Toshimichi Ikemura**  [orcid.org/0000-0001-9931-172X](https://orcid.org/0000-0001-9931-172X)

Department of Bioscience, Nagahama Institute of Bio-Science and Technology. Shiga-ken 526-0829, Japan

## REFERENCES

- Abe, T, Akazawa, Y, Toyoda, A, Niki, H and Baba, T.** 2020. Batch-Learning Self-Organizing Map Identifies Horizontal Gene Transfer Candidates and Their Origins in Entire Genomes. *Frontiers in microbiology*, 11: 1486. DOI: <https://doi.org/10.3389/fmicb.2020.01486>
- Abe, T, Kanaya, S, Kinouchi, M, Ichiba, Y, Kozuki, T and Ikemura, T.** 2003. Informatics for unveiling hidden genome signatures. *Genome research*, 13: 693–702. DOI: <https://doi.org/10.1101/gr.634603>
- Benvenuto, D, Giovanetti, M, Salemi, M, Prosperi, M, De Flora, C, Junior Alcantara, LC, Angeletti, S and Ciccozzi, M.** 2020. The global spread of 2019-nCoV: a molecular evolutionary analysis. *Pathogens and global health*, 114: 64–67. DOI: <https://doi.org/10.1080/20477724.2020.1725339>

- Elbe, S** and **Buckland-Merrett, G.** 2017. Data, disease and diplomacy: GISAID's innovative contribution to global health. *Global challenges (Hoboken, NJ)*, 1: 33–46. DOI: <https://doi.org/10.1002/gch2.1018>
- Gorbalenya, AE, Baker, SC, Baric, RS, de Groot, RJ, Drosten, C, Gulyaeva, AA, Haagmans, BL, Lauber, C, Leontovich, AM, Neuman, BW, et al.** 2020. The species Severe acute respiratory syndrome-related coronavirus: classifying 2019-nCoV and naming it SARS-CoV-2. *Nature microbiology*, 5: 536–544. DOI: <https://doi.org/10.1038/s41564-020-0695-z>
- Hadfield, J, Megill, C, Bell, SM, Huddleston, J, Potter, B, Callender, C, Sagulenko, P, Bedford, T and Neher, RA.** 2018. Nextstrain: real-time tracking of pathogen evolution. *Bioinformatics (Oxford, England)*, 34: 4121–4123. DOI: <https://doi.org/10.1093/bioinformatics/bty407>
- Han, AX, Parker, E, Scholer, F, Maurer-Stroh, S and Russell, CA.** 2019. Phylogenetic Clustering by Linear Integer Programming (PhyCLIP). *Molecular Biology and Evolution*, 36: 1580–1595. DOI: <https://doi.org/10.1093/molbev/msz053>
- Hu, B, Guo, H, Zhou, P and Shi, ZL.** 2021. Characteristics of SARS-CoV-2 and COVID-19. *Nature reviews Microbiology*, 19: 141–154. DOI: <https://doi.org/10.1038/s41579-020-00459-7>
- Ikemura, T, Iwasaki, Y, Wada, K, Wada, Y and Abe, T.** 2021. AI for the collective analysis of a massive number of genome sequences: various examples from the small genome of pandemic SARS-CoV-2 to the human genome. *Genes Genet Syst. (in press)*. DOI: <https://doi.org/10.1101/2021.05.23.445371>
- Ikemura, T, Wada, K, Wada, Y, Iwasaki, Y and Abe, T.** 2020. Unsupervised explainable AI for simultaneous molecular evolutionary study of forty thousand SARS-CoV-2 genomes. *bioRxiv*. DOI: <https://doi.org/10.1101/2020.10.11.335406>
- Iwasaki, Y, Abe, T and Ikemura, T.** 2021. Human cell-dependent, directional, time-dependent changes in the mono- and oligonucleotide compositions of SARS-CoV-2 genomes. *BMC Microbiol*, 21: 89. DOI: <https://doi.org/10.1186/s12866-021-02158-6>
- Iwasaki, Y, Abe, T, Wada, K, Itoh, M and Ikemura, T.** 2011. Prediction of directional changes of influenza A virus genome sequences with emphasis on pandemic H1N1/09 as a model case. *DNA research*, 18: 125–136. DOI: <https://doi.org/10.1093/dnares/dsr005>
- Iwasaki, Y, Abe, T, Wada, Y, Wada, K and Ikemura, T.** 2013. Novel bioinformatics strategies for prediction of directional sequence changes in influenza virus genomes and for surveillance of potentially hazardous strains. *BMC infectious diseases*, 13: 386. DOI: <https://doi.org/10.1186/1471-2334-13-386>
- Kanaya, S, Kinouchi, M, Abe, T, Kudo, Y, Yamada, Y, Nishi, T, Mori, H and Ikemura, T.** 2001. Analysis of codon usage diversity of bacterial genes with a self-organizing map (SOM): characterization of horizontally transferred genes with emphasis on the E. coli O157 genome. *Gene*, 276: 89–99. DOI: [https://doi.org/10.1016/S0378-1119\(01\)00673-4](https://doi.org/10.1016/S0378-1119(01)00673-4)
- Karlin, S, Campbell, AM and Mrzsek, J.** 1998. Comparative DNA analysis across diverse genomes. *Annual review of genetics*, 32: 185–225. DOI: <https://doi.org/10.1146/annurev.genet.32.1.185>
- Kirby, T.** 2021. New variant of SARS-CoV-2 in UK causes surge of COVID-19. *The Lancet Respiratory medicine*, 9: e20–e21. DOI: [https://doi.org/10.1016/S2213-2600\(21\)00005-9](https://doi.org/10.1016/S2213-2600(21)00005-9)
- Kohonen, T.** 1990. The self-organizing map. *Proceedings of the IEEE*, 78: 1464–1480. DOI: <https://doi.org/10.1109/5.58325>
- Kohonen, T, Oja, E, Simula, O, Visa, A and Kangas, J.** 1996. Engineering applications of the self-organizing map. *Proceedings of the IEEE*, 84: 1358–1384. DOI: <https://doi.org/10.1109/5.537105>
- Kumar, S, Stecher, G, Li, M, Knyaz, C and Tamura, K.** 2018. MEGA X: Molecular Evolutionary Genetics Analysis across Computing Platforms. *Mol Biol Evol*, 35: 1547–1549. DOI: <https://doi.org/10.1093/molbev/msy096>
- Mangeat, B, Turelli, P, Caron, G, Friedli, M, Perrin, L and Trono, D.** 2003. Broad antiretroviral defence by human APOBEC3G through lethal editing of nascent reverse transcripts. *Nature*, 424: 99–103. DOI: <https://doi.org/10.1038/nature01709>
- Mercatelli, D and Giorgi, FM.** 2020. Geographic and Genomic Distribution of SARS-CoV-2 Mutations. *Frontiers in microbiology*, 11: 1800. DOI: <https://doi.org/10.3389/fmicb.2020.01800>
- Qiu, Y, Abe, T, Nakao, R, Satoh, K and Sugimoto, C.** 2019. Viral population analysis of the taiga tick, *Ixodes persulcatus*, by using Batch Learning Self-Organizing Maps and BLAST search. *The Journal of veterinary medical science*, 81: 401–410. DOI: <https://doi.org/10.1292/jvms.18-0483>
- Simmonds, P.** 2020. Rampant C→U Hypermutation in the Genomes of SARS-CoV-2 and Other Coronaviruses: Causes and Consequences for Their Short- and Long-Term Evolutionary Trajectories. *mSphere*, 5. DOI: <https://doi.org/10.1128/mSphere.00408-20>
- Sun, J, He, WT, Wang, L, Lai, A, Ji, X, Zhai, X, Li, G, Suchard, MA, Tian, J, Zhou, J, et al.** 2020. COVID-19: Epidemiology, Evolution, and Cross-Disciplinary Perspectives. *Trends in molecular medicine*, 26: 483–495. DOI: <https://doi.org/10.1016/j.molmed.2020.02.008>
- Tang, X, Wu, C, Li, X, Song, Y, Yao, X, Wu, X, Duan, Y, Zhang, H, Wang, Y, Qian, Z, et al.** 2020. On the origin and continuing evolution of SARS-CoV-2. *National Science Review*, 7: 1012–1023. DOI: <https://doi.org/10.1093/nsr/nwaa036>
- Wada, K, Wada, Y and Ikemura, T.** 2020. Time-series analyses of directional sequence changes in SARS-CoV-2 genomes and an efficient search method for candidates for advantageous mutations for growth in human cells. *Gene*: X, 5: 100038. DOI: <https://doi.org/10.1016/j.gene.2020.100038>

**Wang, R, Chen, J, Gao, K, Hozumi, Y, Yin, C and Wei, GW.** 2021. Analysis of SARS-CoV-2 mutations in the United States suggests presence of four substrains and novel variants. *Communications biology*, 4: 228. DOI: <https://doi.org/10.1038/s42003-021-01754-6>

**World Health Organization.** 2020. Coronavirus Disease (COVID-2019). *Situation Reports*. URL: <https://www.who.int/emergencies/diseases/novel-coronavirus-2019>.

Abe et al.  
*Data Science Journal*  
DOI: 10.5334/dsj-2021-029

12

**TO CITE THIS ARTICLE:**

Abe, T, Furukawa, R, Iwasaki, Y, Ikemura, T. 2021. Time-Series Trend of Pandemic SARS-CoV-2 Variants Visualized Using Batch-Learning Self-Organizing Map for Oligonucleotide Compositions. *Data Science Journal*, 20: 29, pp. 1–12. DOI: <https://doi.org/10.5334/dsj-2021-029>

Submitted: 23 March 2021

Accepted: 10 September 2021

Published: 21 September 2021

**COPYRIGHT:**

© 2021 The Author(s). This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International License (CC-BY 4.0), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited. See <http://creativecommons.org/licenses/by/4.0/>.

*Data Science Journal* is a peer-reviewed open access journal published by Ubiquity Press.

