

1 Running title: SARS-CoV-2 genome clusters analyzed by Deep Learning

2

## 3 **Cluster Analysis of SARS-CoV-2 Gene using Deep Learning** 4 **Autoencoder: Gene Profiling for Mutations and Transitions**

5

6

7 Jun Miyake<sup>\*1,2,3</sup>, Takaaki Sato<sup>1,2</sup>, Shunsuke Baba<sup>1,2</sup>, Hayao Nakamura<sup>1,2</sup>, Hirohiko Niioka<sup>4</sup>,

8

Yoshihisa Nakazawa<sup>2</sup>

9

<sup>1</sup> Department of Material and Life Science, Graduate School of Engineering, Osaka University,

10

<sup>2</sup> Hitz Research Alliance Laboratory, Graduate School of Engineering, Osaka University

11

<sup>3</sup> Osaka University Global Center for Medical Engineering and Informatics

12

<sup>4</sup> Osaka University Institute for Dataability Science

13

14 \* Corresponding Author (email-address: [jun\\_miyake@bpe.es.osaka-u.ac.jp](mailto:jun_miyake@bpe.es.osaka-u.ac.jp))

15

Keywords: Autoencoder, Deep Learning, SARS-CoV-2, Genome, Mutation, Classification, Cluster,

16

17

### 18 **Abstract**

19

We report on a method for analyzing the variant of coronavirus genes using autoencoder.

20

Since coronaviruses have mutated rapidly and generated a large number of genotypes,

21

an appropriate method for understanding the entire population is required. The method

22

using autoencoder meets this requirement and is suitable for understanding how and

23

when the variants emerge and disappear. For the over 30,000 SARS-CoV-2 ORF1ab gene

24

sequences sampled globally from December 2019 to February 2021, we were able to

25

represent a summary of their characteristics in a 3D plot and show the expansion,

26

decline, and transformation of the virus types over time and by region. Based on ORF1ab

27

genes, the SARS-CoV-2 viruses were classified into five major types (A, B, C, D, and E in

28

the order of appearance): the virus type that originated in China at the end of 2019 (type

29

A) practically disappeared in June 2020; two virus types (types B and C) have emerged in

30

the United States and Europe since February 2020, and type B has become a global

31

phenomenon. Type C is only prevalent in the U.S. and is suspected to be associated with

32

high mortality, but this type also disappeared at the end of June. Type D is only found in

33

Australia. Currently, the epidemic is dominated by types B and E.

34

## 35 **Introduction**

36 The coronavirus outbreak at the end of 2019 has had unprecedented and significant consequences.  
37 Various researches have been conducted to understand the global trend of genetic alterations (Gómez-  
38 Carballa et al. 2020; Jones and Manrique 2020; Nie et al. 2020; Rochman et al. 2020). Technologies  
39 for the analysis of viruses and other genomes have been developed mainly in the field of molecular  
40 biology for basic research. High-speed gene sequencing technology has enabled the analysis of more  
41 than 40,000 cases worldwide in one year (NCBI Nucleotide Database; NCBI Virus Database). In  
42 order to understand the alternation of viral genomes while utilizing the huge amount of information,  
43 it would be helpful to conceptualize these viral mutations and visualize the spatiotemporal transition.

44 We have been studying the application of deep-learning autoencoder for analyzing gene  
45 sequences (Miyake et al. 2018). The feature extraction capability of autoencoder is useful for this  
46 kind of analysis. There is no need to organize the potentially characteristic sites in the gene  
47 beforehand. In our previous study of the human leukocyte antigen A (HLA-A) gene, we discovered  
48 that autoencoder can correctly represent and classify differences in HLA-A alleles (Miyake et al.  
49 2018). Autoencoder has the potential to extract the genetic characteristics of a gene at a level close to  
50 human recognition. A brand-new method of classification could be realized.

51 By using a deep learning autoencoder, various analyses of genes can be performed in a limited  
52 period of time using a GPU computer, as long as the target is about tens of thousands of genes with  
53 the length of a coronavirus genome (tens of thousands of base pairs). Autoencoder does not require a  
54 gene pre-processing, such as alignment and marking of characteristic gene sequences, nor the need  
55 to prepare supervised learning data in advance. Despite this, gene types can be classified and  
56 displayed as clusters in three-dimensional space. Similar genes in sequences form a single cluster and  
57 the group can be intuitively grasped. The spatial distances between genes/clusters can serve as an  
58 indicator of genetic relationships and may contribute to a sophisticated understanding of evolutionary  
59 processes.

60 In this paper, we used the ORF1ab gene sequences of the new coronaviruses (collected  
61 between December 2019 and February 2021), which were obtained from the NCBI Virus and NCBI  
62 Genbank databases, to extract the self-contained features of about 30,000 genes and display them in  
63 three-dimensional space to investigate how the SARS-CoV-2 virus mutated over time.

64

## 65 **Methods**

66 The Tensor-Flow library (V2.0 downgrade to V1.0) was used as Deep Learning for autoencoder. The  
67 computer was constructed in our laboratory equipped with a GPU (NVIDIA Quadro P-6000), i7  
68 CPU, 64GB RAM memory. OS is Windows 10 or Linux (Ubuntu) OS. The learning was usually 1

69 million times. The ORF1ab gene location of each gene in the NCBI nucleotide database (NCBI  
70 Genbank Database) was determined using the reference sequence (NC\_045512.2).

71 In order to achieve high accuracy in the analysis using autoencoder, it is desirable that the  
72 length of the sample genes is uniform within a certain range. The length of the genomes in the  
73 database is difficult to use because the sequencing methods are different and not uniform. The  
74 ORF1ab gene contains the major part of non-structural proteins (15 types) and occupies more than  
75 about 2/3 of the entire viral genome. The distribution of the length of the genome or the large number  
76 of undecided sequences makes the scattered dots wider and the 3D plot harder to see. Genome  
77 samples with two or more consecutive undetermined RNA sequences were excluded from the  
78 analysis.

79 The nucleotide sequence data of the new coronaviruses were obtained from the NCBI Virus  
80 Database (NCBI Virus Database). The sequencing was downloaded on February 26, 2021 (the sample  
81 collection dates correspond to December 19, 2019 to February 16, 2021.). Data that were described only  
82 up to the year and data that did not specify the collection site were excluded from the analysis. Data with  
83 collection dates only up to the month were assigned a day in the middle of the month; for the two genes with  
84 no date found in December 2019, we assigned December 19. Finally, 31050 ORF1ab gene was extracted  
85 and analyzed for characteristics (cluster analysis) and time series.

86 Genetic analysis by autoencoder that we have already reported (Miyake et al. 2018) was used  
87 in this study. Namely, we applied the document vector method (the nucleotide sequence was replaced  
88 by a vector ( $4^5 = 1,024$  dimensions) with a normalized histogram of 1,024 words consisting of 5-mer  
89 tiny nucleotide sequences without alignment). In this research, the hierarchy was compressed to four  
90 layers and three dimensions. In order to visualize the obtained 3D data, we plotted them as x, y, z  
91 coordinates in 3D space. Each dot corresponds to a variant nucleotide sequence. The spatial distance  
92 from the center of the all dots plotted in 3D space was calculated and used to represent the gene  
93 profile in a time series.

94 Phylogenetic trees were constructed using maximum likelihood phylogenetic analysis  
95 (RAxML) with 1000 bootstraps (GENETYX ver. 15, GENETYX Co., Tokyo, Japan). Alignment of  
96 nucleotide sequences was performed using the above software.

97

## 98 **Results**

99 The ORF1ab genes, extracted from the genomes of 33,915 novel coronaviruses (12/19/2019–  
100 02/16/2021), were categorized into eight clusters in 3D space (Fig. 1). The variation of the ORF1ab  
101 gene sequence length was small, leading to the result that the separation of the clusters was clear. The  
102 3D-compressed dots correspond to respective RNA strands as many as the number of samples used.

103           These dots are plotted at different spatial locations, but instead of being simply distributed,  
104 several types of sets (clusters) appear. The ORF1ab genes were shown to form the eight clusters with  
105 similar characteristics in 3D coordinates and distances from the center. Close proximity of three pairs  
106 of neighboring clusters suggested their similarities in mutation profiles, respectively. The eight  
107 clusters of the ORF1ab genes were categorized into five major groups (Fig. 1). These clusters were  
108 named A, B, C, D, and E in the order of appearance.

109           In order to investigate the temporal changes, we replotted the 3D dots monthly or bimonthly  
110 for the collection period (Fig. 2). The ORF1ab genes collected during the two months of December  
111 2019 and January 2020 showed a predominance of type A cluster in the center of 3D plotted genes  
112 (Fig. 2 **a, b**). Type B became the dominant genotype from February to March, and type E became the  
113 dominant genotype from April to May. Type C started in February and fell and disappeared in June.  
114 Type D appeared in June-July and disappeared in October.

115           The time series of the type C obtained by autoencoder analysis seems to be consistent with  
116 the emergence and disappearance and geo location of coevolving variant group 4 (CEVg4) reported  
117 by Chan et al. (Chan et al. 2020). Based on genome frequencies and geo locations, our classification  
118 of types A1, A2, B1, and D seemed to correspond to the wild type, CEVg3, CEVg1, and CEVg6,  
119 respectively. The B2 cluster is in a different location from the B1 cluster and, is a group of similar  
120 size to the B1 cluster (Figs. 1 and 2). In contrast, there is no CEVg similar to CEVg1.

121           The distance of each dot from the center of the all dots in the 3D space was calculated and  
122 used to represent the genotype profile in a time series by country/region. The data were color-coded  
123 by cluster and displayed separately by geographic region (Figs. 3 and 4).

124           The stretching and extinction of genotypes was quite frequent, with a new species emerging  
125 and disappearing approximately every two months. It is unclear whether this was derived from a  
126 single species, or whether a species that originally existed was grown.

127           The maximum likelihood phylogenetic trees of 88 ORF1ab genes and their corresponding  
128 full-length genomes are shown in Fig. 5**a** and **b**, respectively. Genes corresponding to genotypes A1,  
129 A2, B1, B2, C, D, E1, and E2 in both the ORF1ab and full-length genomes of the phylogenetic tree  
130 are represented by the same colors as in Fig. 4. Genes classified into clusters can be regarded as  
131 having a certain degree of correlation. This is a fairly good correlation considering the fact that they  
132 have different principles.

133

## 134 **Discussion**

135           The deep-learning autoencoder was able to successfully classify the genotypes of SARS-CoV-2  
136 viruses, and the autoencoder method is useful for overarching classification, which is similar to  
137 human cognitive abilities. It is widely recognized that the genes of coronaviruses change one after

138 another, and the autoencoder method is a useful method for easily recognizing time-series changes.  
139 As shown in Figs. 2–4, it is a simple and straightforward method that allows us to grasp the elongation  
140 and disappearance of clusters in the viral genome.

141 Judging from the present analysis, the occurrence and disappearance of new species appears  
142 to be observed about every two months. Such correlations are available for understanding: type A  
143 first appeared in December, 2019, but largely disappeared by the end of June, 2020. The most  
144 widespread strains globally appear to be types B and E. Type B also appears to have started close to  
145 type A and spread, possibly as a result of successive changes in each.

146 Coronaviruses are RNA viruses and are particularly rapidly-mutated genomes. As shown in  
147 the eight clades in Fig. 1, mutations do not cause simple spread, but lead to the formation of clusters.  
148 Mutant species that pop out of the clusters form new clusters there as well. We can read a form of  
149 repeated expression and flourishing of new species in nature.

150 The cluster classification by the autoencoder method (shown in Fig. 1) showed a certain  
151 correlation with the classification by the phylogenetic tree method (Fig. 5). In both ORF1ab and the  
152 whole genome, a certain degree of cohesion was observed for genotypes A1, A2, B1, B2, C, D, E1,  
153 and E2, and we judged that there is considerable correlation in gene sequencing. Because both  
154 classification by autoencoder and phylogenetic tree analysis based on sequence homology and  
155 differences, the methods are do not always match perfectly in principle, but they help each other to  
156 understand classification. As shown in Fig. 5, they can be considered as essentially distant  
157 correlations as classification methods for gene sequences.

158 We found the eight clusters using over 30,000 SARS-CoV-2 ORF1ab genes in the NCBI  
159 Virus database, whereas Chan et al. identified nine CEVg using 86,450 genomes in the GISAID  
160 database (Chan et al. 2020). Yet there is not enough data to rigorously compare the differences  
161 between the ABCDE and CEVg classifications. autoencoder-based classification is considered to be  
162 a useful method for scanning the entire SARS-CoV-2 virus for variations or for rapid genetic  
163 classification of viral genes and viruses with a certain genetic distance.

164 With regard to the new coronavirus, more than 40,000 gene sequencings were performed in  
165 one year for the whole world. In order to take advantage of the vast information space made possible  
166 by next-generation sequencing technology, we believe that we need technology to grasp the entire  
167 picture of genetic variation and its distribution patterns. We hope that artificial intelligence will  
168 contribute to the development of methods for rapid recognition and classification of genetic  
169 mutations. We believe that being able to explain the direction of mutations and the principles that  
170 constrain them will make a significant contribution to this field. A better understanding of viral  
171 evolution will allow us to respond more effectively and quickly to pandemics.

172 We are currently working on a detailed analysis of the internal structure of the autoencoder  
173 cluster and would like to point out that there may be new applications for classification.

174

## 175 **Author contributions**

176 JM designed the project; JM and TS wrote the manuscript; SB developed computer system and  
177 software; HM collected data; HN supervised the study of artificial intelligence; YN realized the  
178 project and scientifically supervised the work.

179

## 180 **Acknowledgements**

181 We should like to express our thanks to Yuta Nitada, Keita Fukuda and Takuto Shimazaki of Osaka  
182 University for programming and computer operations. This work was supported partially by Global  
183 Center for Medical Engineering and Informatics of Osaka University, Hitz Research Alliance  
184 Laborator of Osaka University, and Japan Agency for Medical Research and Development (Grant  
185 Number 20bm0804008h0004.PI. Prof. S. Miyagawa of Medical School of Osaka University).

186

## 187 **References**

- 188 Chan AP, Choi Y, Schork NJ. Conserved genomic terminals of SARS-CoV-2 as coevolving functional  
189 elements and potential therapeutic targets. *mSphere* 5: e00754-20 (2020). doi: 10.1128/mSphere.00754-20.
- 190 Coronaviridae Study Group of the International Committee on Taxonomy of Viruses. The species Severe acute  
191 respiratory syndrome-related coronavirus: classifying 2019-nCoV and naming it SARS-CoV-2. *Nat.*  
192 *Microbiol.* (2020). doi: 10.1038/s41564-020-0695-z.
- 193 Gómez-Carballa A, Bello X, Pardo-Seco J, Martínón-Torres F and Salas A. Mapping genome variation of  
194 SARS-CoV-2 worldwide highlights the impact of COVID-19 super-spreaders. *Genome Res.* 30:1434–1448  
195 (2020). .doi: 10.1101/gr.266221.120.
- 196 Jones LR, Manrique JM. Quantitative phylogenomic evidence reveals a spatially structured SARS-CoV-2  
197 diversity. *Virology* 550:70-77 (2020). doi: 10.1016/j.virol.2020.08.010. Epub 2020 Aug 26.
- 198 Miyake J, Kaneshita Y, Asatani S, Tagawa S, Niioka H, Hirano T. Graphical classification of DNA sequences  
199 of HLA alleles by Deep learning. *Human Cell* 31:102–105 (2018). doi: 10.1007/s13577-017-0194-6.
- 200 NCBI Genbank Database: <https://www.ncbi.nlm.nih.gov/genbank/>
- 201 NCBI Virus Database: <https://www.ncbi.nlm.nih.gov/labs/virus/vssi/#/find-data/virus>
- 202 Nie Q, Li X, Chen W, Liu D, Chen Y, Li H, Li D, Tian M, Tan W, Zai J. Phylogenetic and analyses of SARS-  
203 CoV-2. *Virus Res.* 287:198098 (2020). doi: 10.1016/j.virusres.2020.198098. Epub 2020 Jul  
204 17.phylodynamic
- 205 Rochman ND, Wolf YI, Faure G, Zhang F, Koonin EV. Ongoing adaptive evolution and globalization of Sars-  
206 Cov-2. *bioRxiv.* 2020 Oct 13:2020.10.12.336644 (2020). doi: 10.1101/2020.10.12.336644.

207 Zhang T, Wu Q, Zhang Z. Probable Pangolin Origin of SARS-CoV-2 Associated with the COVID-19  
208 Outbreak. *Curr Biol.* 30:1346-1351.e2 (2020). doi: 10.1016/j.cub.2020.03.022.  
209 Zhou P, Yang XL, Wang XG, Hu B, Zhang L, Zhang W, Si HR, Zhu Y, Li B, Huang CL, Chen HD, Chen J,  
210 Luo Y, Guo H, Jiang RD, Liu MQ, Chen Y, Shen XR, Wang X, Zheng XS, Zhao K, Chen QJ, Deng F, Liu  
211 LL, Yan B, Zhan FX, Wang YY, Xiao GF, Shi ZL. A pneumonia outbreak associated with a new coronavirus  
212 of probable bat origin. *Nature* (2020), doi: 10.1038/s41586-020-2012-7.  
213

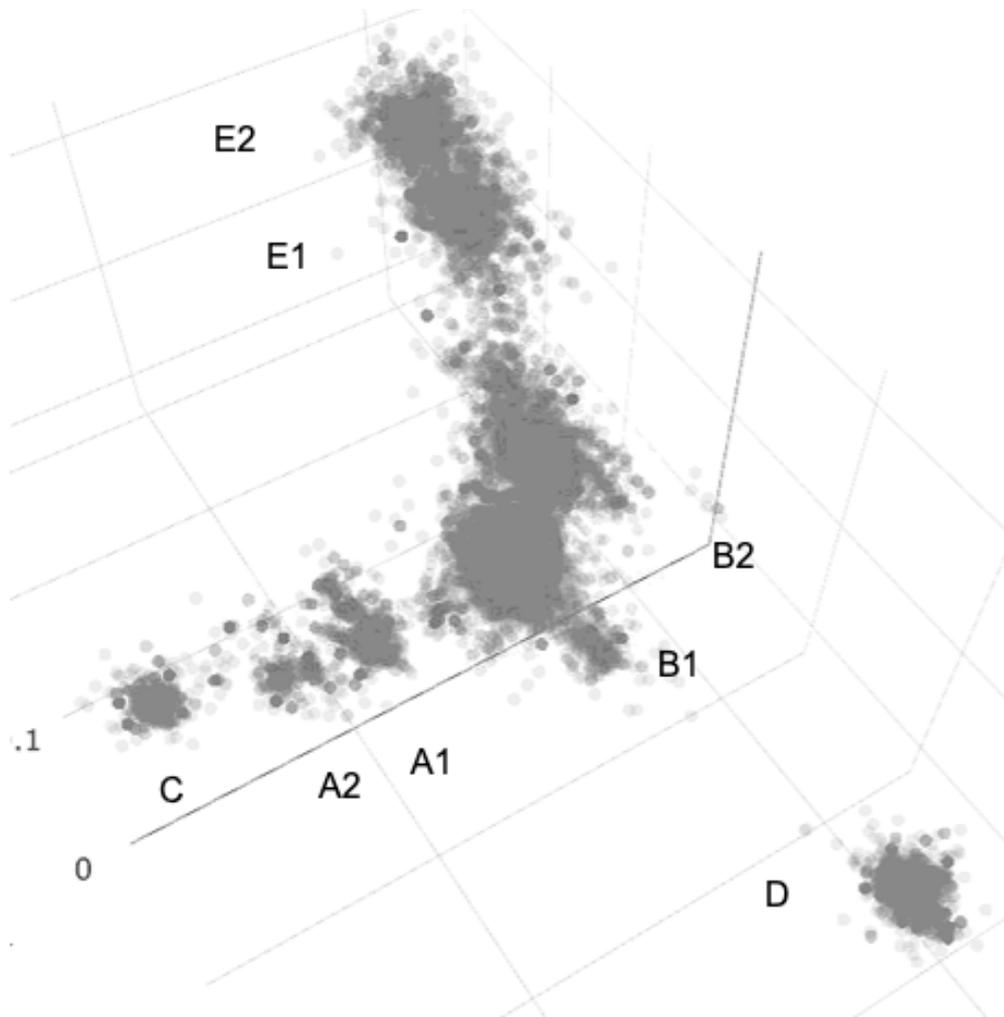
214 **Figures**

215

216

217

218

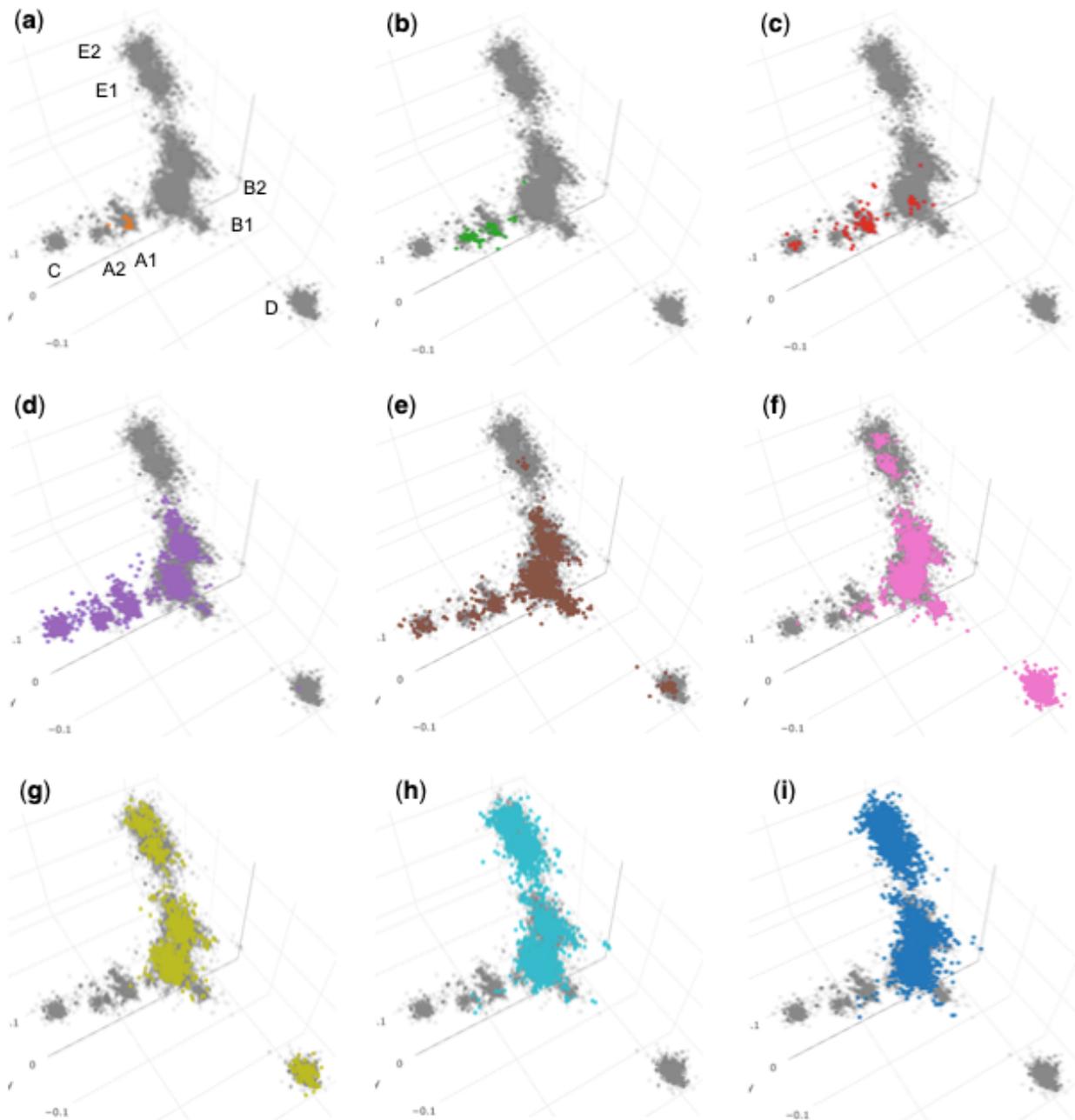


219

220

221 **Fig. 1. Three-dimensional plot of the ORF1ab genes.**

222 A deep-learning autoencoder classified 33,915 ORF1ab genes of SARS-CoV-2 viruses into eight  
223 clusters. The symbols of the clusters are given in the order of the time of emergence. Two clusters  
224 that appeared at the same time differed with suffix. ORF1ab genes were dissected from 33,915  
225 genome sequences of SARS-CoV-2 viruses collected from December 2019 to February 2021.  
226 Occurrence time was shown by colored dots monthly (December 2019–February 2020) or  
227 bimonthly (March 2020–February 2021).



228

229 **Fig. 2. Monthly or bimonthly trend of gene clusters.**

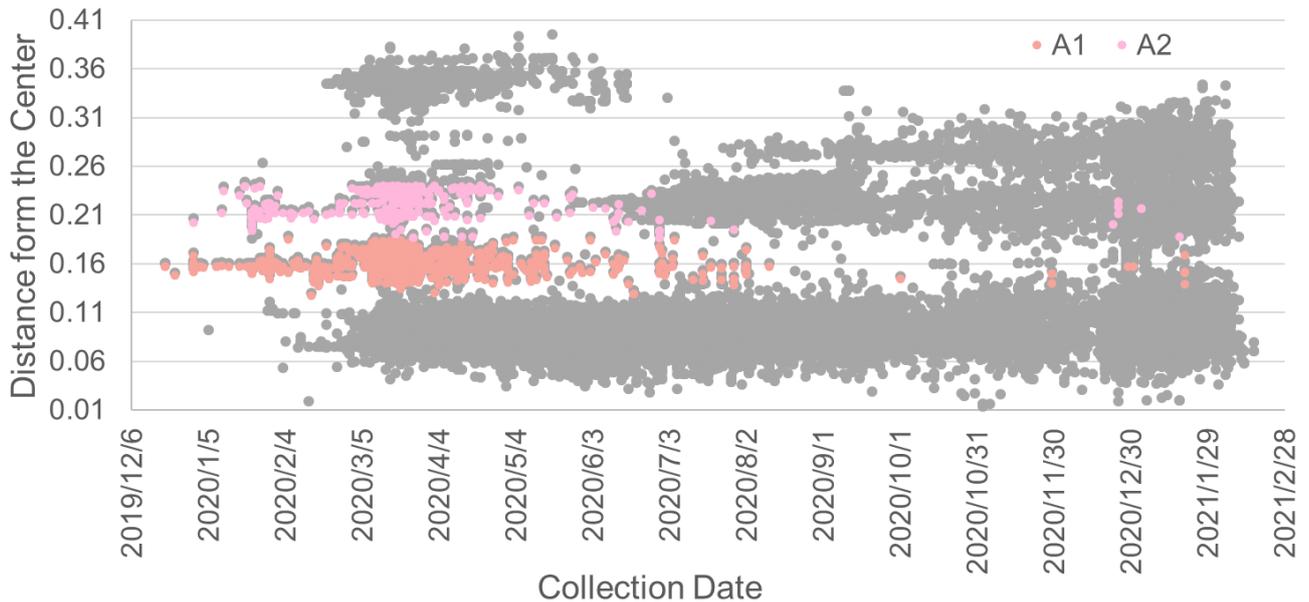
230 (a) ORF1ab genes of SARS-CoV-2 viruses collected in December 2019 (orange). (b) January 2020  
231 (green). (c) February 2020 (red). (d) March–April 2020 (purple). (e) May–June 2020 (brown). (f)  
232 July–August 2020 (pink). (g) September–October 2020 (yellow green). (h) November–December  
233 2020 (cyan). (i) January–February 2021 (blue). Shown as background in light gray is the ORF1ab  
234 gene for the entire period.

235

236

237

238



239

240

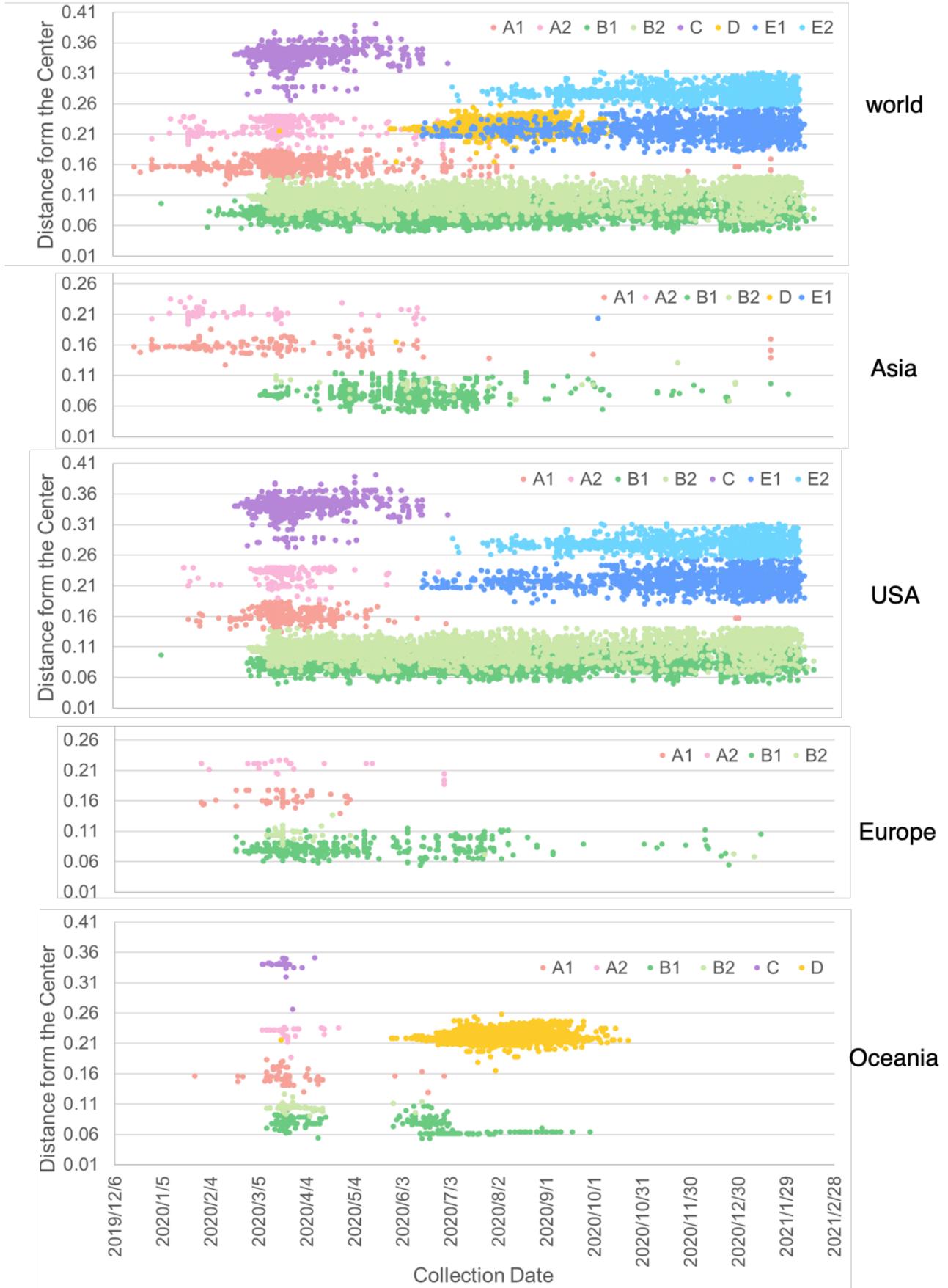
241 **Fig. 3. Time course plot of the SARS-CoV-2 ORF1ab gene for A1 and A2**

242 **clusters.**

243 This figure shows the way of visualization of temporal transitions in target gene clusters. The gray  
244 or colored dots, which represent the ORF1ab genes classified by autoencoder, were plotted by the  
245 collection date and distance from the center in the 3D plot in Fig. 1. Based on the coordinate data of  
246 each gene, the time evolution of a particular cluster can be shown by coloring. Here we show an  
247 example of A1 and A2 extracted.

248

249



250

251

252 **Fig. 4. Emergence and transition of clusters.**

253 Autoencoder-classified ORF1ab genes for World, Asia, USA, Europe and Oceania are separately  
254 illustrated. Clusters are extracted as shown an example in Fig. 2. Dot colors represent different  
255 clusters (A1, red; A2, pink; B1, green; B2, light green; C, purple; D, orange; E1, blue; E2, light  
256 blue). To prevent mixing, dots in border regions between the clusters were omitted. Along the axis  
257 of the distance from the center, B2 cluster had some overlap with B1 cluster.



260 **Fig. 5. Phylogenetic trees of SARS-CoV-2 ORF1ab genes and whole genomes.**

261 (a) Phylogeny of 88 ORF1ab genes of SARS-CoV-2 viruses. (b) Phylogeny of 88 whole genomes  
262 of the SARS-CoV-2 viruses. Eleven genes from each of eight clusters were randomly selected from  
263 the central region of the respective clusters of A1, A2, B1, B2, C, D, E1 and E2 in the autoencoder  
264 3D plot (Fig.1). The accession numbers of the viruses used were identical between the ORF1ab  
265 genes and whole genomes. Phylogenies showed similar features between ORF1ab genes and whole  
266 genomes that members from single or two clusters formed isolated or mixed clades. Phylogenetic  
267 trees were constructed with maximum likelihood method with RAxML, 1,000 bootstraps  
268 (GENETYX ver. 15).