

TCS-TR-B-10-43

TCS Technical Report

Master's Thesis: Virus Data Clustering based on
Kolmogorov Complexity

by

YU ZHU

Division of Computer Science

Report Series B

December 22, 2010



Hokkaido University
Graduate School of
Information Science and Technology

Email: zhuyu07@ist.hokudai.ac.jp

Phone: +81-011-706-7675

Fax: +81-011-706-7675

Summary

Influenza viruses are probably a major cause of morbidity and mortality world wide. Large segments of the human population are affected every year. In June 2009, World Health Organization declared the influenza due to a new strain of swine origin H1N1 was responsible for the 2009 influenza pandemic. And on June 11, the WHO declared an H1N1 pandemic moving the alert level to phase 6, marking the first global pandemic since the 1968 Hong Kong influenza. There are a lot of data mining methods used in biological sciences to analysis viruses. But if one designs data mining algorithms based on domain knowledge, then the resulting algorithms tend to have many parameters. Determining how relevant particular features are is often difficult and may require a certain amount of guessing.

In this thesis, we introduce a universal data mining method which we call *parameter-free data mining*. The approach of *parameter-free data mining* is aimed at scenarios where we are not interested in a certain similarity measure but in the similarity between objects themselves. The most promising approach to this paradigm is called *normalized information distance* which uses Kolmogorov complexity theory as its basis. As the *normalized information distance* (NID) cannot be computed, we apply this idea to standard compression algorithms, such as `gzip` and `bzip`, have been used as approximations of the Kolmogorov complexity. This yields the *normalized compression distance* (NCD) as approximation of the NID.

To demonstrate the usefulness of the normalized compression distance for clustering influenza viruses data, two kinds of compressors and two clustering algorithms have been used, which verified that this approach neither depend on the compression methods nor the clustering methods we choose.

Chapter 1

Introduction

1.1 Background of Parameter Free Data Mining

Influenza viruses were probably responsible for the disease described by Hippocrates in 412BC [13], and thus they have been with us for a long time. Influenza viruses remain a major cause of morbidity and mortality world wide. To analyze or predict the viruses data, many methods to classify the viruses data have been developed.

The similarity between objects is a fundamental notion in everyday life. It is also fundamental to many data mining and machine learning algorithms, and, in particular to clustering algorithms. Often the similarity between objects is measured by a domain-specific distance measure based on features of the objects. For example, the distance between pieces of music can be measured by using features like rhythm, pitch, or melody, i.e., features that do not make sense in any other domain. To develop such methods one needs special knowledge about the application domain for extracting the relevant features beforehand. Such an approach does not only cause difficulties, but includes a certain danger or risk of being biased.

If one is pursuing the approach to design data mining algorithms based on domain knowledge, then the resulting algorithms tend to have many parameters. By using these parameters, one can then control the algorithms' sensitivity to certain features. Determining how relevant particular features are is often difficult and may require a certain amount of guessing. Expressing this differently, one has to tune the algorithms which is requiring domain knowledge and a larger amount of experience. Furthermore, it may be expensive, error prone and time consuming to arrive at a suitable tuning.

1.2 Main Idea of Kolmogorov Complexity

However, as a radically different approach, the paradigm of parameter-free data mining has emerged (cf. Keogh *et al.* [14]). The main idea of parameter-free data mining is the design of algorithms that have no parameters and that are universally applicable in all areas.

The problem is whether or not such an approach can be realized at all. It is only natural to ask how an algorithm can perform well if it is not based on extracting the important features of the data and if we are not allowed to adjust its parameters until it is doing the right thing. As expressed by Vitányi *et al.* [23], *if we a priori know the features, how to extract them, and how to combine them into exactly the distance measure we want, we should do just that. For example, if we have a list of cars with their color, motor rating, etc. and want to cluster them by color, we can easily do that in a straightforward way.*

So the approach of parameter-free data mining is aiming at scenarios where we are not interested in a certain similarity measure but in the similarity between the objects themselves. The most promising approach to this paradigm uses Kolmogorov complexity theory [16] as its basis.

The key ingredient to this approach is the so-called *normalized information distance* (NID) which was developed by various researchers during the past decade in a series of steps (cf., e.g., [4, 15, 6]).

More formally the *normalized information distance* between two strings x and y is defined as

$$NID(x, y) = \frac{\max\{K(x|y), K(y|x)\}}{\max\{K(x), K(y)\}}, \quad (1.1)$$

where $K(x|y)$ is the length of the shortest program that outputs x on input y , and $K(x)$ is the length of the shortest program that outputs x on the empty input. It is beyond the scope of the present paper to discuss the technical details of the definition of the NID. We refer the reader to Vitányi *et al.* [23].

The NID has nice theoretical properties, the most important of which is universality. The NID is called *universal*, since it accounts for the dominant difference between two objects (cf. Li *et al.* [15] and Vitányi *et al.* [23] and the references therein).

In a sense, the NID captures all computational ways in which the features needed in the traditional approach could be defined. Since its definition involves the Kolmogorov complexity $K(\cdot)$, the NID cannot be computed. Therefore, to apply this idea to real-world data mining tasks, standard compression algorithms, such as `gzip`, `bzip`, or PPMZ, have been used as approximations of the Kolmogorov complexity. This yields the *normalized compression distance* (NCD) as approximation of the NID (cf. Definition 1).

In a typical data mining scenario we are given some objects as input. The pairwise NCDs for all objects in question form a distance matrix. This matrix can be processed further until finally standard algorithms, e. g., clustering algorithms can be applied. This has been done in a variety of typical data mining scenarios with remarkable success. Works of literature and music have been clustered according to genre or author; evolutionary trees of mammals have been derived from their mitochondrial genome; language trees have been derived from several linguistic corpora (cf., e.g., [6, 14, 5, 8, 3]).

As far as virus data are concerned, Cilibrasi and Vitányi [9] used the SARS TOR2 draft genome assembly 120403 from Canada’s Michael Smith Genome Sciences Centre and compared it to other viruses by using the NCD. They used the `bzip` compressor and applied their quartet tree heuristic for hierarchical clustering. The resulting ternary tree showed relations very similar to those shown in the definitive tree based on medical-macrobiological genomics analysis which was obtained later (see [9] for details).

1.3 Procedure

The main goal of the present paper is a detailed analysis of the general method outlined above in the domain of influenza viruses. More specifically, we are interested in learning whether or not specific gene data for the hemagglutinin of influenza viruses are *correctly* classifiable by using the concept of the NCD. For this purpose we have chosen a set of 106 gene sequences from the National Center for Biotechnology Information for which the correct classification of the hemagglutinin is known. As explained in Chapter 3, there are 16 subtypes commonly called H1, . . . , H16. For these 106 gene sequences (or subsets thereof) we then compute the NCD by using the CompLearn Toolkit (cf. [7]) as done in [9].

This computation returns a symmetric matrix D such that d_{ij} is the NCD between the data entries i and j (henceforth called distance matrix). Furthermore, we study the influence of the compressor chosen and restrict ourselves here to the `zlib` and `bzip` compressors which are the standard two built-in compressors for the CompLearn Toolkit.

The next step is the clustering. Here of course the variety of possible algorithms is large. Note that the CompLearn Toolkit contains also an implementation of quartet tree heuristic for hierarchical clustering. However, this heuristic is computationally quite expensive and does currently not allow to handle a matrix of dimension 106×106 . Therefore, we have decided to try the *hierarchical clustering* algorithm from the R package (called `hclust`) with the average option. In this way we obtain a rooted tree showing the relations among the input data.

The second clustering algorithm used is *spectral clustering* via `kLines` (cf. Fischer and Poland [10]). We have successfully applied this method before (cf. [21, 20]) in settings where the NID is approximated by the so-called Google Distance. It should be noted that spectral clustering generally requires the transformation of the distance matrix into an adjacency matrix of pairwise similarities (henceforth called similarity matrix). The clustering is then done by analyzing its spectrum.

The results obtained are generally very promising. Quite often, we obtained a *perfect* clustering independently of the method used. On the other hand, when including all data or a rather large subset thereof, the clustering obtained is not perfect but the number of errors made is still sufficiently small to make the results interesting.

Without going into details here, it can be said that the `zlib` compressor seems more suitable in this setting than the `bzip` compressor (see Section 3.2 for details).

Chapter 2

Background and Theory

2.1 Definitions and Axioms for *normalized compression distance*

As explained in the Introduction, the theoretical basis for computing the distance matrix is deeply based in Kolmogorov complexity theory. In the following we assume the definition of the NID as shown in Equation (1.1). The definition of the NID depends on the function K which is *uncomputable*. Thus, the NID is *uncomputable*, too.

Using a real-word compressor, one can approximate the NID by the NCD (cf. Definition 1). Again, we omit details and refer the reader to [23].

Definition 1 *The normalized compression distance between two strings x and y is defined as*

$$NCD(x, y) = \frac{C(xy) - \min\{C(x), C(y)\}}{\max\{C(x), C(y)\}},$$

where C is any given data compressor.

Common data compressors are `gzip`, `bzip`, `zlib`, etc. Note that the compressor C has to be computable and *normal* in order to make the NCD a useful approximation. This can be stated as follows.

Definition 2 *A compressor C is said to be normal if it satisfies the following axioms for all strings x, y, z and the empty string λ .*

- (1) $C(xx) = C(x)$ and $C(\lambda) = 0$; (identity)
- (2) $C(xy) \geq C(x)$; (monotonicity)
- (3) $C(xy) = C(yx)$; (symmetry)
- (4) $C(xy) + C(z) \leq C(xz) + C(yz)$; (distributivity)

up to an additive $O(\log n)$ term, with n the maximal binary length of a string involved in the (in)equality concerned.

These axioms are in various degrees satisfied by good real-world compressors like `bzip`, `PPMZ` and `gzip`, where the latter did not perform so well, as informal experiments have shown (cf. [6]). Also note that in all cases the compressor-specific window or block size determines the maximum usable length of the arguments. As a matter of fact, for our data these axioms seem to be fulfilled.

Here, we take 3 random viruses sequences from the data set. To verify the compressors like `bzip` and `gzip` are satisfied with the axioms, we compressed viruses sequences respectively, and then compare the compression sizes.

We chose H1N1AF091309 (x), H1N1D10477(y), H1N1U47310(z) from the data set. After compressed, the size of the file(s) are as

	x	$x \cdot x$	$x \cdot y$	$y \cdot x$	z	$x \cdot z + y \cdot z$
<code>gzip</code>	615B	635B	662B	757B	362B	1279B
<code>bzip</code>	569B	590B	816B	806B	341B	1524B

Compression sizes by `gzip` and `bzip`

Based on the compression results, we can say both `bzip` and `zlib` are approximately satisfied with those Axioms. They are *normal* compressors. So for our investigations we used the built-in compressors `bzip` and `zlib` and the `ncd` function from the `CompLearn` Toolkit (cf. [7]). After having done this step, we have a distance matrix $D = (d^{\text{ncd}}(x, y))_{x, y \in X}$, where $X = (x_1, \dots, x_n)$ is the relevant data list.

2.2 Clustering Algorithms

2.2.1 Hierarchical Clustering Algorithm

Next, we turn our attention to clustering. First, we shortly outline the hierarchical clustering as provided by the R package, i.e., by the program `hclust` (cf. [2]). Input is the $(n \times n)$ distance matrix D . The program uses a measure of dissimilarity for the objects to be clustered. Initially, each object is assigned to its own cluster and the program proceeds iteratively. In each iteration the two most similar clusters are joint, and the process is repeated until only a single cluster is left. Furthermore, in every iteration the distances between clusters are recomputed by using the Lance–Williams dissimilarity update formula for the particular method used.

The methods differ in the way in which the distances between clusters are recomputed. Provided are the *complete linkage method*, the *single linkage method*, and the *average linkage clustering*. In the first case, the distance between any two clusters is equal to the greatest similarity from any member of one cluster to any member of the

other cluster. This method works well for compact clusters but causes sensitivity to outliers. The second method pays attention solely to the area where the two clusters come closest to one another. The more distant parts of the clusters and the overall structure of the clusters is not taken into account. If the total number of clusters is large, a messy clustering may result.

The *average linkage clustering* defines the distance between any two clusters to be the average of distances between all pairs of objects from any member of one cluster to any member of the other cluster. As a result, the average pairwise distance within the newly formed cluster, is minimum.

Heuristically, the average linkage clustering should give the best results in our setting, and thus we have chosen it (see also Manning *et al.* [17] for a thorough exposition). Note that for hierarchical clustering the number k of clusters does *not* to be known in advance.

2.2.2 Spectral Clustering Algorithm

Next, the spectral *spectral clustering* algorithm used is shortly explained. Spectral clustering is an increasingly popular method for analyzing and clustering data by using only the matrix of pairwise similarities. It was invented more than 30 years ago for partitioning graphs (cf., e.g., Spielman and Teng [22] for a brief history and Luxburg [24] for a tutorial). Formally, spectral clustering can be related to approximating the normalized min-cut of the graph defined by the adjacency matrix of pairwise similarities [26]. Finding the exactly minimizing cut is an NP-hard problem.

The transformation of the distance matrix into a similarity matrix is done by using a suitable kernel function. In our experiments we have used the Gaussian kernel function, i.e.,

$$k(x, y) = \left(\exp\left(-\frac{1}{2}d(x, y)^2/(2 \cdot \sigma^2)\right) \right), \quad (2.1)$$

where σ is the kernel width. As pointed out by Perona and Freeman [19], there is nothing magical with this function. Moreover, it is most commonly used. An advantage of using the Gaussian kernel function is that the resulting similarity matrix is positive definite.

So, the remaining problem is a suitable choice for σ . Unfortunately, the performance of spectral clustering heavily depends on this σ . In the experiments, we compute the mean value of the entries of the distance matrix D and then set $\sigma = \text{mean}(D)/\sqrt{2}$. In this way, the kernel is most sensitive around $\text{mean}(D)$. Though we are not aware of a theoretical result supporting this choice, it worked remarkably well and further studies are needed to explore the properties of this choice.

The final spectral clustering algorithm for a known number of clusters k is stated below. Following Fisher and Poland [10], we did not use *k-means* here, details had been explained in [10].

Algorithm Spectral clustering of a data list

Input: data list $X = (x_1, x_2, \dots, x_n)$, number of clusters k

Output: clustering $c \in \{1 \dots k\}^n$

1. for $x, y \in X$, compute the distance matrix $D = (d^{\text{ncd}}(x, y))_{x, y \in X}$
2. compute $\sigma = \text{mean}(D)/\sqrt{2}$
3. compute the similarity matrix $A = (\exp(-\frac{1}{2}d(x, y)^2/(2 \cdot \sigma^2)))$
4. compute the Laplacian $L = S^{-\frac{1}{2}}AS^{-\frac{1}{2}}$, where $S_{ii} = \sum_j A_{ij}$ and $S_{ij} = 0$ for $i \neq j$
5. compute top k eigenvectors $V \in R^{n \times k}$
6. cluster V using kLines [10]

Chapter 3

Experiments and Results

In this section we describe the data used, the experiments performed and the results obtained.

3.1 Influenza Viruses - The Data Set

3.1.1 Brief Introduction of Influenza Viruses

We shortly describe the data set used. For any relevant background concerning the biological aspects of the influenza viruses we refer the reader to Palese and Shaw [18] and Wright *et al.* [25].

Influenza viruses were probably a major cause of morbidity and mortality world wide. Large segments of the human population are affected every year. The family of *Orthomyxoviridae* is defined by viruses that have a negative-sense, single-stranded, and segmented RNA genome. There are five different genera in the family of *Orthomyxoviridae*: the influenza viruses A, B and C; *Thogotovirus*; and *Isavirus*. Influenza A viruses have a complex structure and possess a lipid membrane derived from the host cell (cf. Figure 3.1).

Biologists classify influenza A viruses primarily by their hemagglutinin (HA) subtypes and neuraminidase (NA) subtypes. So far, 16 subtypes of HA are known and commonly denoted by H1, . . . , H16. In addition to these HA types, biologists distinguish 9 NA subtypes denoted by N1, . . . , N9.

Influenza A viruses of all 16 hemagglutinin (H1-H16) and 9 neuraminidase (N1-N9) subtypes are maintained in their nature host, i.e., the duck. Of these duck viruses, H1N1, H2N2 and H3N2 subtypes jumped into human population, and caused three pandemics in the last century. Therefore, in the experiments performed we have exclusively selected data of influenza viruses that have been obtained from viruses hosted by the duck.

The complete genome of these influenza viruses has 8 segmented-genes. Of these 8 genes, here we are only interested in their HA gene, since HA is a the major target

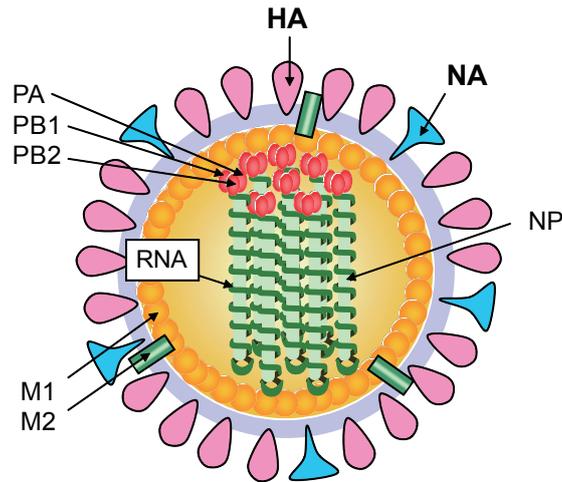


Figure 3.1: Influenza A virus

of antibodies that neutralize viral infectivity, and responsible for binding the virus to the cell it infects. The corresponding gene is found on segment 4.

Each datum consists of a sequence of roughly 1800 letters from the alphabet $\{A, T, G, C\}$, e.g., looking such as *AAAAGCAGGGGAATTCACAAT TAAACAAAAT...TGTATATAATTAGCAA*. These gene sequences are publicly available from the National Center for Biotechnology Information (NCBI) which has one of the largest collections of such sequences (cf. [11]).

When analyzed by biologists the definite method to determine the correct HA subtype is based on the antiserum that prevent the docking of the virus. Sometimes biologists also compare the actual sequence to already analyzed sequences and produce a guess based on the Hamming distance of the new sequence to the analyzed ones.

3.1.2 Data set

As explained in the Introduction, the primary goal of the investigations undertaken is to cluster the sequences correctly with respect to their HA subtype. In order to achieve this goal with collected from each subtype up to 8 examples. The reason for choosing at most 8 sequences from each type has been caused by their availability. While for some subtypes there are many sequences, there are also subtypes for which only very few sequences are available. The extreme case is the subtype H16 for which only one sequence is in the data base. Table 3.2 shows the number of sequences chosen. All of the subtypes from H1–H7 have 8 sequences for each, but some of the subtypes like H8, H13, H14 and H15, they do not have enough sequences, even H16, there is only one sequence available. So the reason we chose 8 sequences for most of the subtypes is to keep the balance between 16 subtypes. And we also did other experiments to prove that increase the sequences number will not help us to get better results.

H1	H2	H3	H4	H5	H6	H7	H8
8	8	8	8	8	8	8	7
H9	H10	H11	H12	H13	H14	H15	H16
8	8	8	8	2	4	4	1

Figure 3.2: Number of sequences for each subtype

It should be noted that most of these sequences are marked as `complete cds`, but some are also marked as `partial cds` by the NCBI. For a complete list of the data description we refer the reader to the dataset in the Appendix.

For the ease of presentation, in the following we use the following abbreviation for the data entries. Instead of giving the full description, e.g.,

```
>gi|113531192|gb|AB271117| /Avian/4 (HA)/H10N1/Hong Kong/1980/// Influenza
A virus (A/duck/Hong Kong/938/80(H10N1)) HA gene for haemagglutinin, complete
cds.
```

we refer to this datum as H10N1AB271117 for short.

Among the available files have been two containing only a very short partial sequence of the gene, i.e., H7N1AM157391 and H10N4AM922160 (483 and 80 letters, respectively). So, we did not consider these two files, since they do not seem to contain enough information.

3.2 Results

3.2.1 Experiments Enviroment

All experiments have been performed under SuSE Linux. As already mentioned, for the hierarchical clustering we used the open source R package (cf. [2]).

The Algorithm (Spectral clustering of a data list) has been realized by performing Step 1 via the `CompLearn` function `ncd` (cf. [7]). Steps 2 through 6 have been implemented in `GNU Octave`, version 2.1.72 (cf. [1]). It should be noted that `ncd` assigns 0.000000 to all elements on the main diagonal of the distance matrix (Version 1.1.5).

3.2.2 Aim

By performing our experiments we aimed to answer the following questions. First, does the NCD provide enough information to obtain a correct clustering for the virus data? Second, does the rather large number of clusters (recall that we 16 HA types) cause any problems? Third, do the answers to the fist and second question depend on the compressor and clustering, respectively, chosen?

3.2.3 Clustering HA sequences for H1 through H3

To get started and for the sake of comparison, we used the subset containing all data belonging to H1, H2, and H3, i.e., a total of 24 sequences (cf. Figure 3.2).

Using the `maketree` program from the CompLearn Toolkit, we get the following clustering (cf. Figures 3.3 and 3.4). As Figures 3.3 and 3.4 show, the data are clearly and correctly separated into three clusters. However, the intra-cluster dissimilarities clearly differ from inter-cluster dissimilarities in Figure 3.3, i.e., for the `zlib` compressor, while there is no such clear difference for the `bzip` compressor (cf. Figure 3.4).

Using `hclust` we obtained the trees shown in Figure 3.5 and 3.6 for the matrix D computed for the compressor `zlib` and `bzip`, respectively.

After having computed the matrix D , we get the following order of the data

H2N4CY003984, H3N1CY005943, H3N2AB277754, H1N9CY017275,
 H1N9CY035248, H3N3CY005936, H2N2L11128, H2N2L11136,
 H2N2L11137, H2N1CY017693, H2N1CY021125, H2N3L11138,
 H1N6CY004458, H1N1D10477, H2N3CY014710, H3N2EU74652,
 H3N2CY006026, H1N1AF091309, H1N1U47310, H3N3AB292410,
 H3N2D21171, H3N2M73771, H1N5CY004498, H1N5CY014968

Since spectral clustering is a hard clustering method, it has to return for each data entry just one class label. Assigning canonically the clusters 1, 2, and 3 to the HA subtypes, we therefore should get the sequence

```

2 3 3 1 1 3 2 2 2 2 2
1 1 2 3 3 1 1 3 3 3 1 1

```

which was indeed returned for both compressors. Note that $\sigma = 0.56078$ and $\sigma = 0.57329$ for the `zlib` and `bzip` compressor, respectively.

3.2.4 Clustering HA sequences for H1 through H8 and H9 through H16

Next, we tried all HA sequences for H1 through H8 and from H9 through H16. The reason for this partition has been caused by the different number of sequences available. Recall that there are only two sequences for H13 and only one sequence for H16 (cf. Figure 3.2).

For H1 through H8 the hierarchical clustering was error free for the `zlib` compressor but not for `bzip` compressor (1 error) (see Figures 3.7 and 3.8).

Interestingly, for H9 through H16 the tree obtained for the `zlib` compressor contains 4 errors, while the one obtained for `bzip` compressor has only one error.

Our spectral clustering algorithm returned a perfect clustering for all HA sequences for H1 through H8 for both compressors. On the other hand, for all sequences from H9 through H16 the results differed with respect to the compressors.

For the `zlib` compressor we obtained 5 errors and for `bzip` the number of errors was 7 when using for σ the mean as described above. However, it is well-known that spectral clustering is quite sensitive to the kernel width σ . So, we also tried to vary it a bit around the mean by rounding it to two decimal digits and then changing the second one. For `zlib` the mean was 0.60873 and after two variations we found $\sigma = 0.59$ which resulted in just one error, i.e., H16 was classified as H13. For the `bzip` compressor such an improvement could not be obtained.

3.2.5 Clustering HA sequences for H1 through H12

As a possible explanation we conjecture that one needs a certain minimum of available sequences in order to arrive at a correct spectral clustering. And there is no clustering method putting one single member into one cluster. Trying all HA sequences for H1 through H12 kind of confirmed this conjecture, since we again obtained a perfect spectral clustering for both compressors.

As Figures 3.9 and 3.10 show, for the hierarchical clustering, the tree obtained for the `zlib` compressor is correct, but the the one obtained for the `bzip` compressor has one error.

3.2.6 Clustering all HA sequences for 16 subtypes

Finally, we tried all data. Again hierarchical clustering was best for the `zlib` compressor and showed only 2 errors. For the `bzip` compressor, we obtained 3 errors (see Figures 3.11 and 3.12).

On the other hand, the best result we could obtain for spectral clustering had 5 errors (for both compressors). Below we show the clustering obtained for the `zlib` compressor for $\sigma = 0.63$, where $c0$ is the desired classification and sp the one returned from the spectral clustering algorithm.

$c0 =$	7	7	14	2	11	12	12	3	7	10
$sp =$	7	7	14	2	11	12	12	3	7	10
$c0 =$	10	5	9	9	9	3	1	1	9	11
$sp =$	10	5	9	9	9	3	1	1	9	11
$c0 =$	11	5	3	7	5	2	2	2	4	10
$sp =$	11	5	3	7	5	2	2	2	4	10
$c0 =$	5	8	12	2	2	4	4	4	11	9
$sp =$	5	8	12	2	2	4	4	4	11	9
$c0 =$	10	2	6	6	6	5	1	1	4	10
$sp =$	10	2	13	13	6	5	1	1	4	10
$c0 =$	7	4	8	15	2	9	9	16	10	14
$sp =$	7	4	8	15	2	9	9	3	10	14
$c0 =$	14	7	7	6	14	7	8	8	12	12
$sp =$	14	7	7	6	14	7	8	8	12	12
$c0 =$	11	15	3	15	5	11	3	1	1	8
$sp =$	11	15	3	15	5	11	3	1	1	8
$c0 =$	4	3	3	6	12	10	4	5	3	6
$sp =$	4	3	3	6	12	10	4	5	3	6
$c0 =$	13	13	12	1	1	11	12	8	11	10
$sp =$	13	13	12	1	1	11	12	8	11	10
$c0 =$	5	9	15	8	6	6				
$sp =$	5	9	15	8	13	13				

So, the errors occur at positions 43, 44, 58, 105, and 106 and affect H6 which is four times assigned to H13 and one time H16 which got in the H3 cluster. We omit further details due to the lack of space.

Note that one can also compute the sum square error (s.s.e.) of all eigenvalues w.r.t. their means in order to determine quite reliably from the eigenvalues of the Laplacian the number k of clusters.

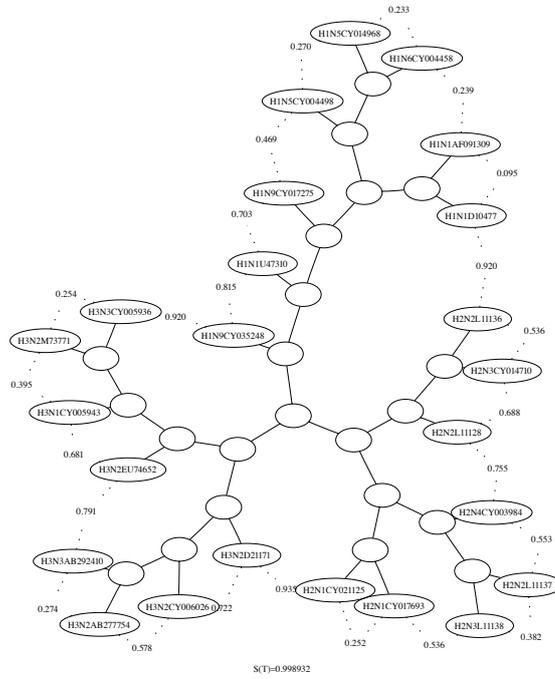


Figure 3.3: Classification of HA sequences for H1 through H3; compr.: zlib

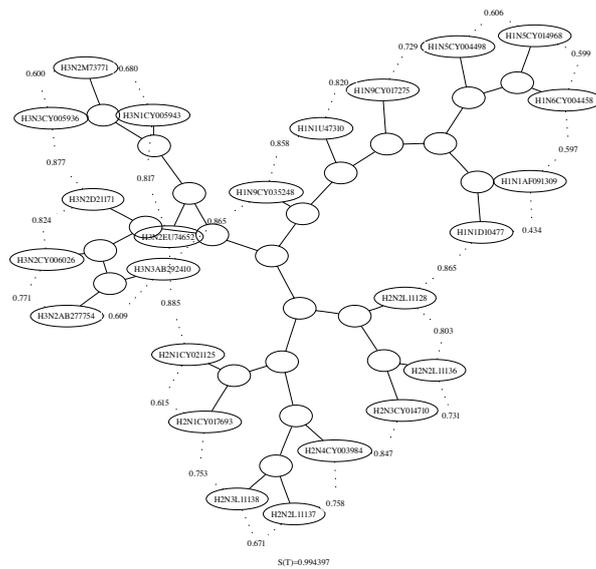


Figure 3.4: Classification of HA sequences for H1 through H3; compr.: bzip

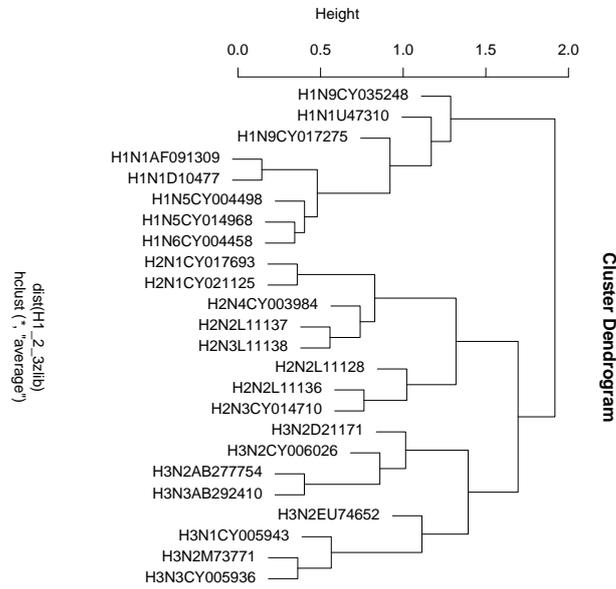


Figure 3.5: Clustering all HA sequences for H1 through H3 via hclust; compr.: zlib

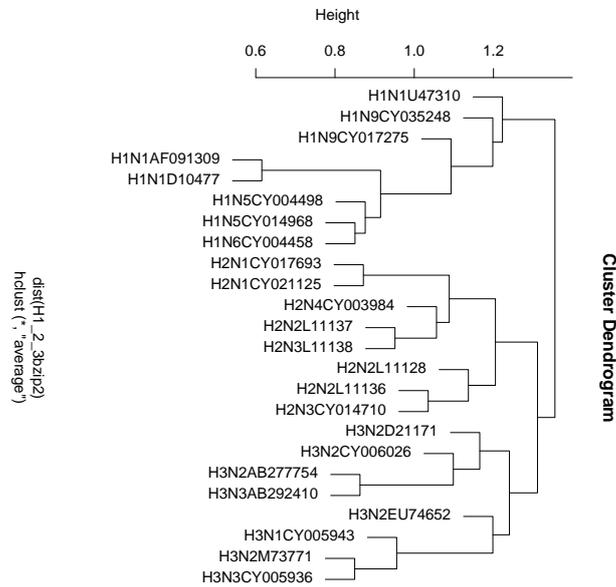


Figure 3.6: Clustering all HA sequences for H1 through H3 via hclust; compr.: bzip

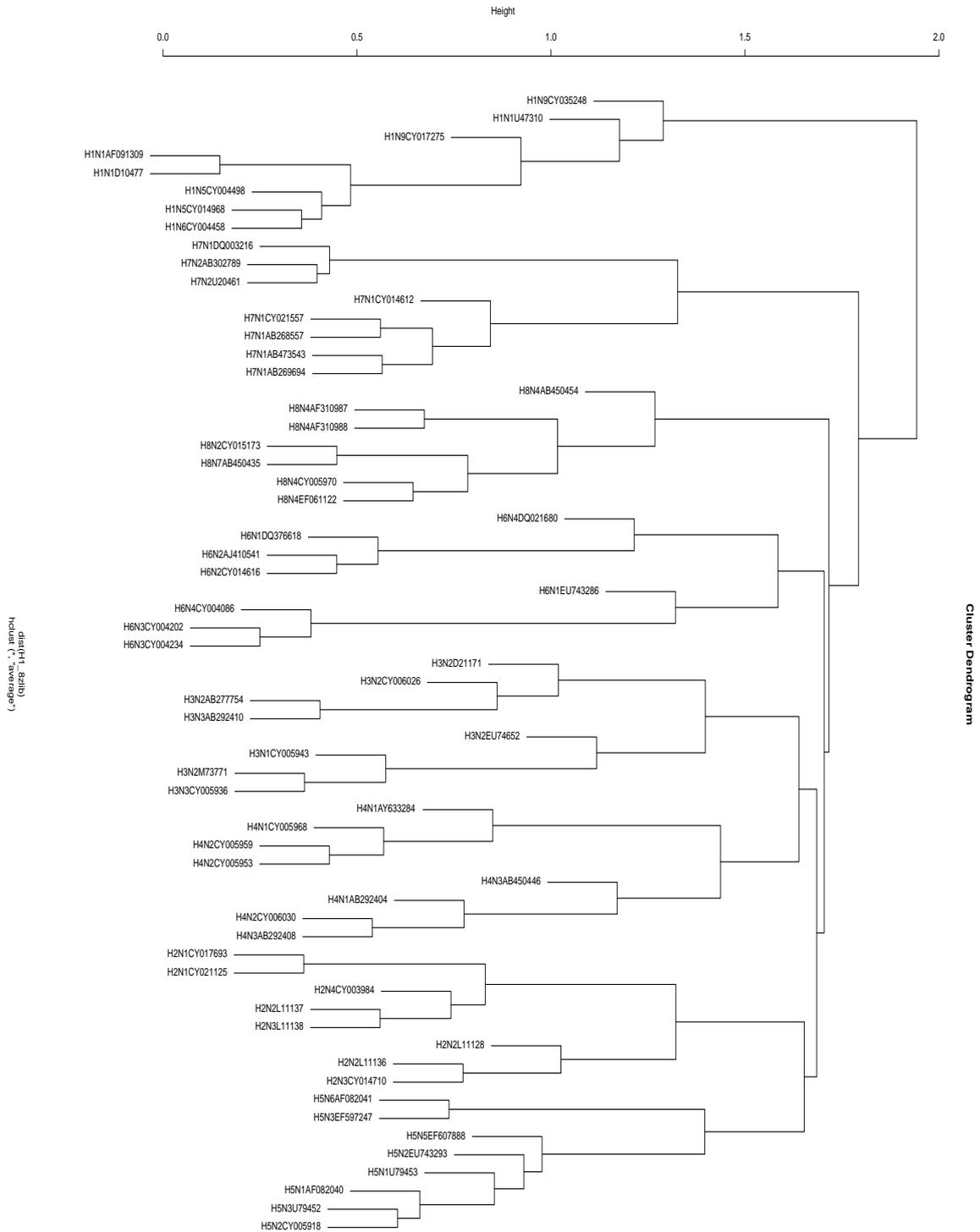


Figure 3.7: Clustering of all HA sequences for H1 through H8 via `hclust`; `compr::zlib`;

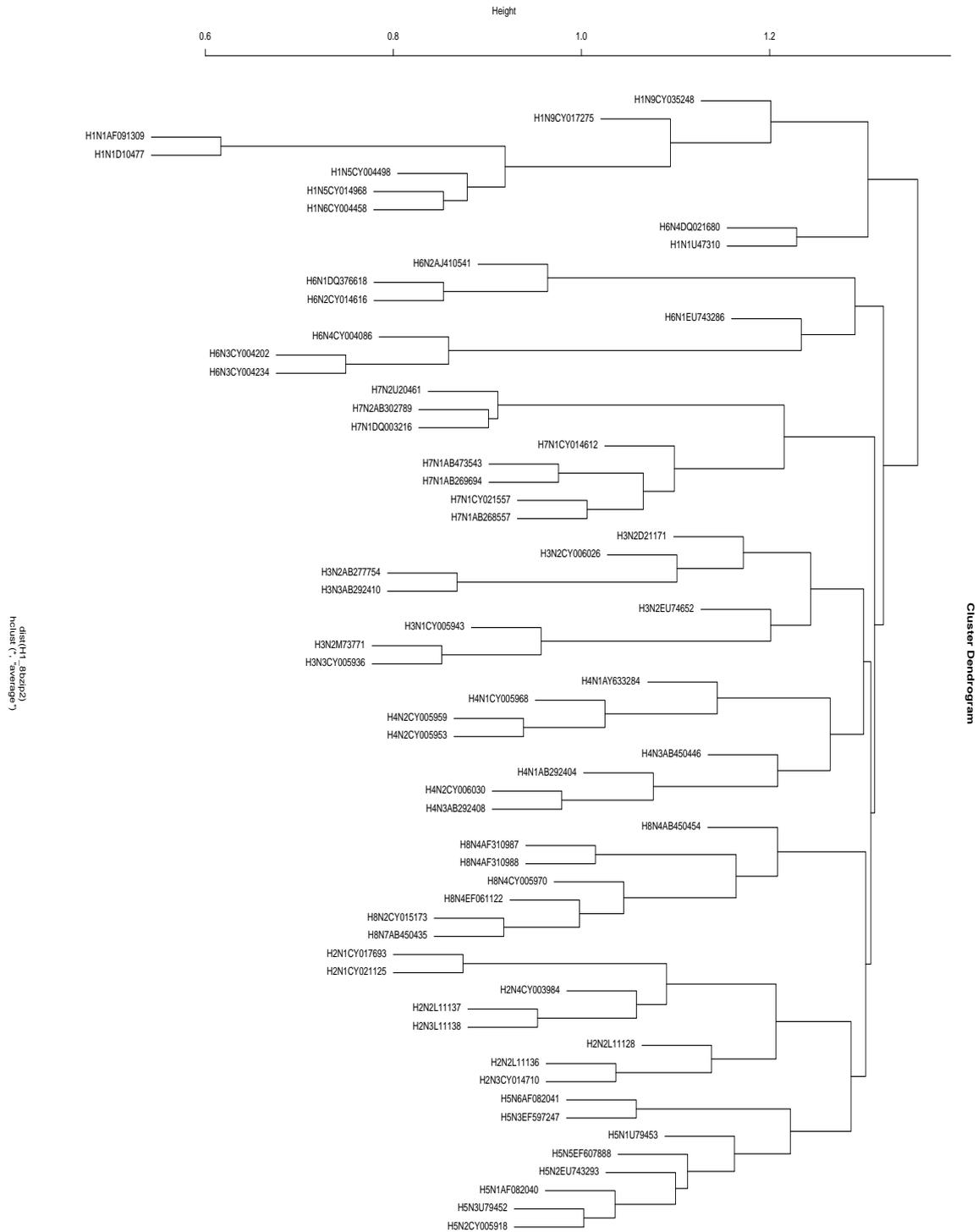


Figure 3.8: Clustering of all HA sequences for H1 through H8 via hclust; compr.:bzip;

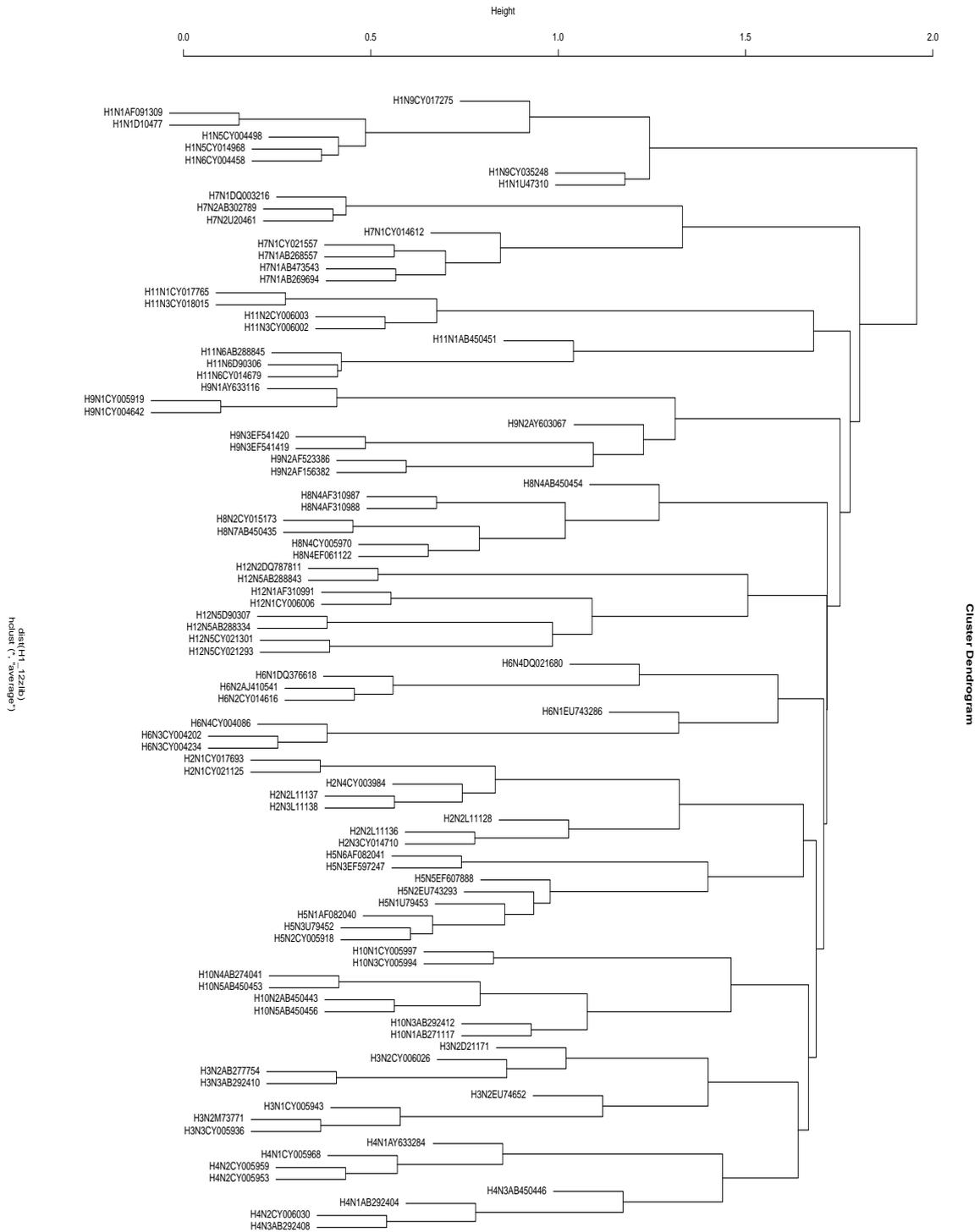


Figure 3.9: Clustering of all HA sequences for H1 through H12 via hclust; compr::zlib;

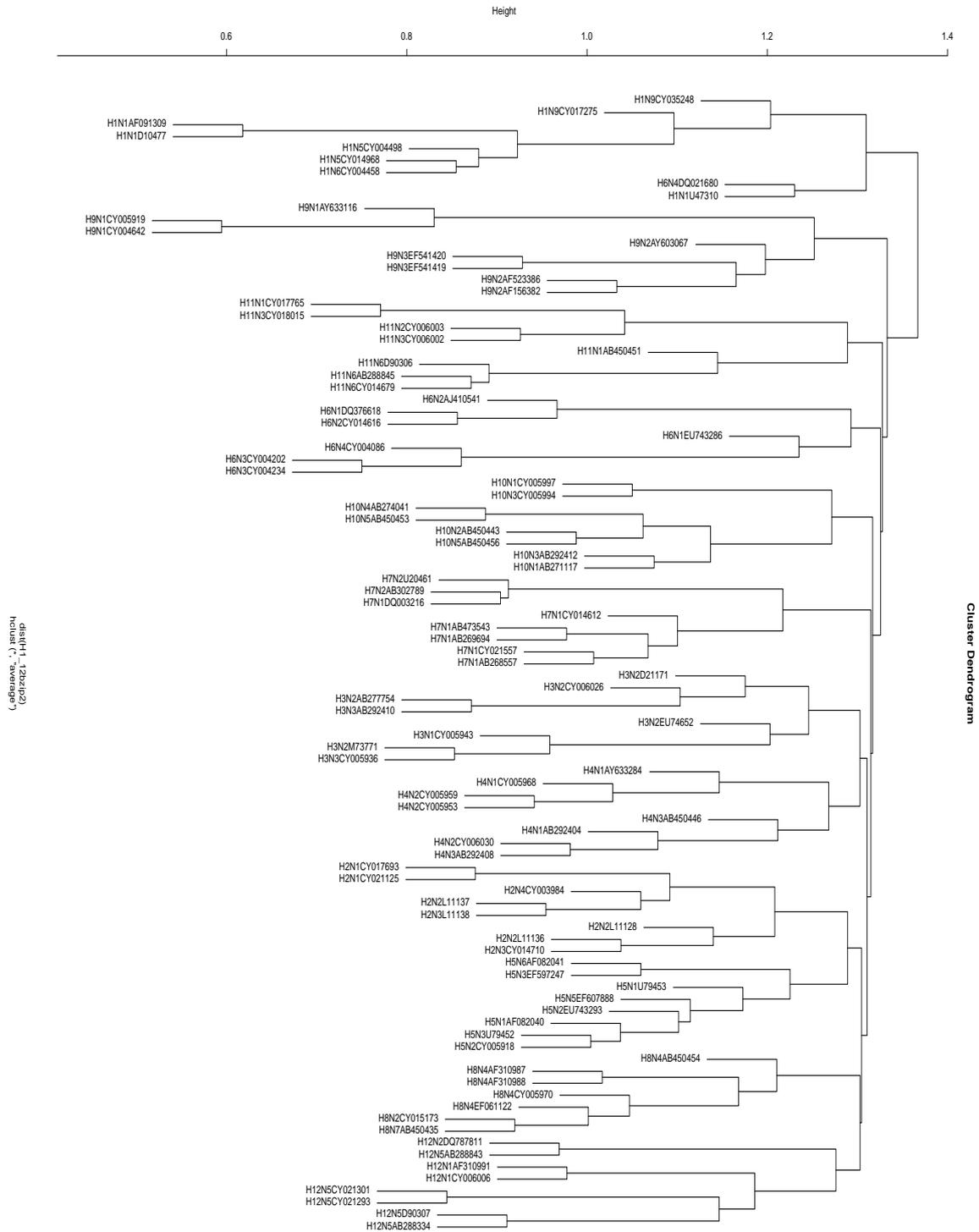


Figure 3.10: Clustering of all HA sequences for H1 through H12 via hclust; compr.:bzip;

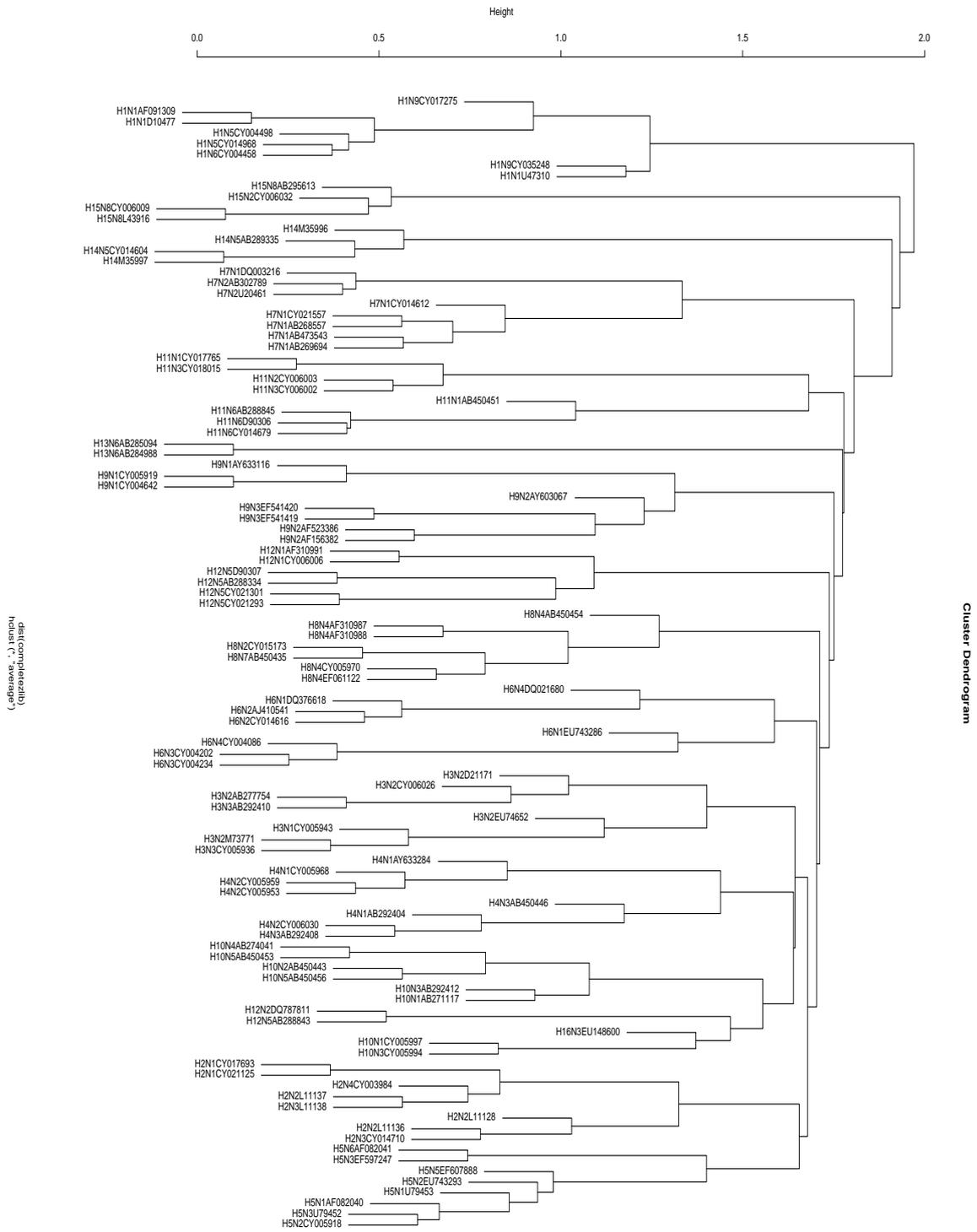


Figure 3.11: Clustering of all HA sequences for all subtypes via `hclust`; `compr::zlib`;

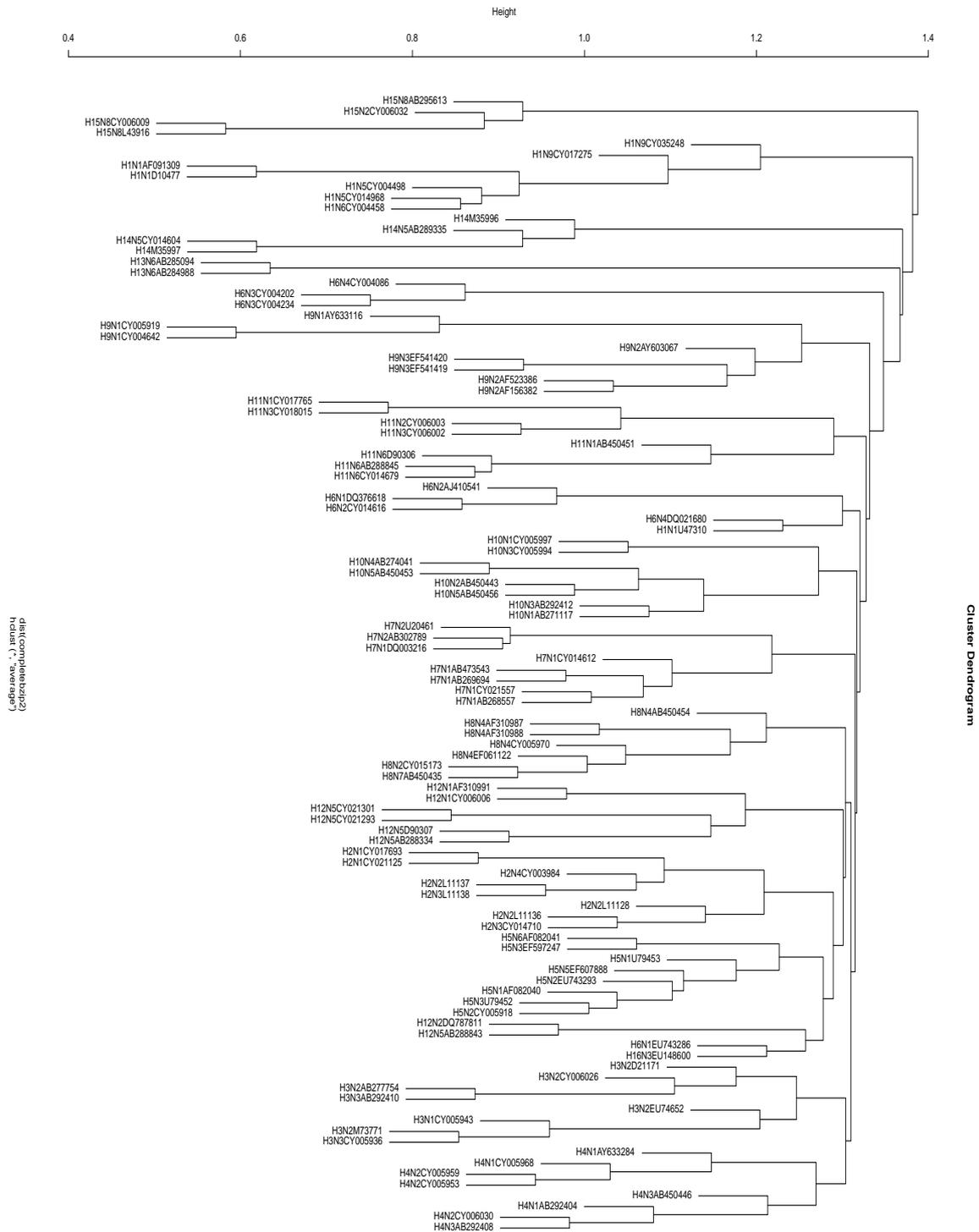


Figure 3.12: Clustering of all HA sequences for all subtypes via hclust; compr.:bzip;

Chapter 4

Application in Swine Influenza Viruses

Our experiments gave us very good clustering results for clustering all 16 subtypes HA sequences. In this Chapter, we do not only focus on the influenza viruses hosted by the nature host **duck** but also other hosts, like human.

Swine influenza virus is common throughout **swine** populations worldwide. Transmission of the virus from **swine** to humans is not common and does not always lead to human influenza, often resulting only in the production of antibodies in the blood. If transmission does cause human influenza, it is called zoonotic **swine** influenza. People with regular exposure to **swine** are at increased risk of **swine** influenza infection. H1N1 is a subtype of influenza A and the most common cause of influenza in humans. Some strains of H1N1 are endemic in humans and cause a small fraction of all influenza-like illness and a large fraction of all seasonal influenza. Other strains of H1N1 are endemic in **swine** which is referred to as **swine** influenza and in birds, namely, avian influenza. But H1N1 strains caused roughly half of all human flu infection in 2006. In June 2009, World Health Organization declared that flu due to a new strain of **swine** origin H1N1 was responsible for the 2009 flu pandemic. This strain is commonly called "swine flu" by the public media. On June 11, 2009, the WHO declared an H1N1 pandemic, moving the alert level to phase 6, marking the first global pandemic since the 1968 Hong Kong flu.

4.1 Latest Global Pandemic—swine influenza

4.1.1 Clustering viruses hosted by human and duck

In this section, we aim to learn which property is more "similar" in these influenza viruses, the *host* or the *subtype*. Usually, the birds can pass avian influenza viruses to **swine**, where the two viruses co-mingle and form a new strain then pass to human or other hosts. But sometimes, the avian influenza viruses or **swine** influenza viruses can also pass to human directly. At this time, the similarity between *subtype* should

be more obvious than the *host*. There is also another possibility for the emergence of new viruses, which we can see from viruses evolution. The new subtype viruses may be the mutation of the subtypes, which already existed in the host. In this case, the similarity between *host* should be more obvious than the *subtype*.

Because of this, we did following experiments. First, we downloaded 8 following **swine** influenza viruses segments 4 from NCBI, and all of them are hosted by *human*.

GQ247726 Influenza A virus (A/Moscow/02/2009(H1N1)) segment 4 hemagglutinin (HA) gene, complete cds

GQ223408 Influenza A virus (A/Beijing/501/2009(H1N1)) segment 4 hemagglutinin (HA) gene, complete cds

GQ221788 Influenza A virus (A/Arizona/02/2009(H1N1)) segment 4 hemagglutinin (HA) gene, complete cds

GQ219577 Influenza A virus (A/Kobe/1/2009(H1N1)) segment 4 hemagglutinin (HA) gene, complete cds

GQ249333 Influenza A virus (A/Paris/2591/2009(H1N1)) segment 4 hemagglutinin (HA) gene, complete cds

GQ219580 Influenza A virus (A/Osaka-C/1/2009(H1N1)) segment 4 hemagglutinin (HA) gene, complete cds

GQ162170 Influenza A virus (A/Mexico/4108/2009(H1N1)) segment 4 hemagglutinin (HA) gene, complete cds

GQ250161 Influenza A virus (A/Guangdong/03/2009(H1N1)) segment 4 hemagglutinin (HA) gene, complete cds

All of the latest **swine** influenza viruses are H1N1 subtype, and because we want to know the relation between *hosts* and between *subtypes*. We reset the data set as H1 subtype viruses and hosted by *human* and **duck**.

Again, first we calculate the symmetric distance matrix, and then use clustering methods to get the cluster results. Again for the ease of presentation, in the following we added a character "d" before the viruses hosted by **duck** and a character "h" before the **swine** influenza viruses because they are hosted by **human**.

As Figure 4.1 and Figure 4.2 show, the viruses data hosted by **human** and **duck** are clearly and correctly separated into two clusters. And the similarity between viruses hosted by human are seem bigger than the viruses hosted by **duck**. That maybe because we chose all viruses hosted by human from 2009, but the viruses hosted by **duck** are different from year by year. However now we are not only satisfied with the methods give us the perfect clustering results, but also want to know whether the results can supply more informations to us.

To answer our question, in influenza viruses which property is more "similar", the *host* or the *subtype*? We need viruses data from both different subtypes and different

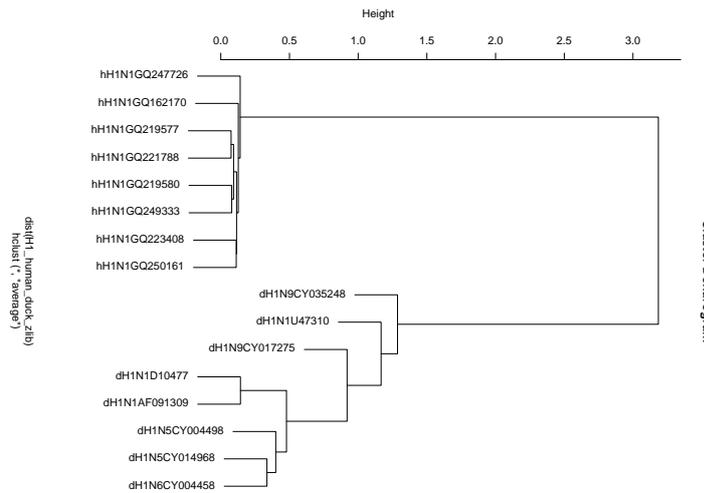


Figure 4.1: Clustering viruses hosted by human and duck via hclust; compr.: zlib

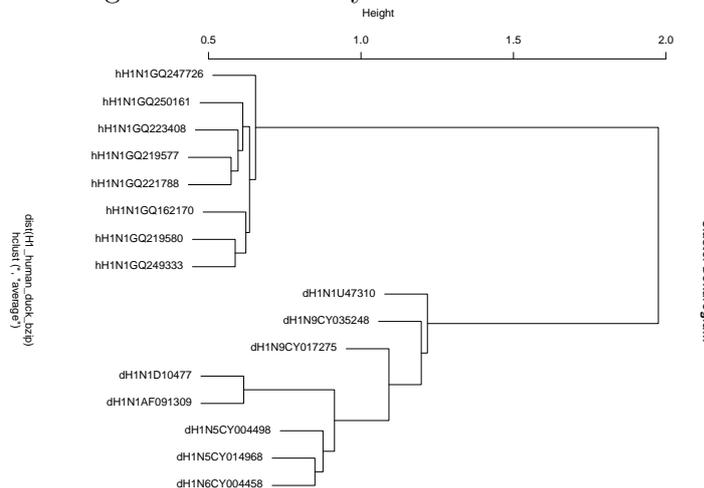


Figure 4.2: Clustering viruses hosted by human and duck via hclust; compr.: bzip

hosts. So we reset the data set again as H1 and H2 subtypes with duck host and human host. Now our data set has 2 kinds of subtypes viruses and 2 kinds of hosts viruses. Then we repeat our experiment procedures. Figure 4.3 and Figure 4.4 give us the results.

Based on the clustering results, we can say, in this data set, it seems like the similarity between subtypes are **bigger** than the similarity between **host**. But we can not say the similarity between subtypes are always **bigger** than **host**. That is also depend on the data we choose. Recall that as WHO declared that influenza due to the new strain of **swine** origin H1N1 was responsible for this pandemic. As we introduced, same subtype viruses maybe jump to different hosts, in this case, it should be more similar in subtype than in host. But if we choose some viruses, they are mutations in genes in the host, or more specifically, in subtype, the similarity in

hosts should be more obviously.

4.2 Clustering viruses with different hosts

As introduced in viruses evolution, **duck** is the original nature host for viruses, how did viruses go to lives in all around the world. One of the guess [12] is when viruses go to **human**, they first go to **swine**. Like the latest pandemic. So we also want to know the relationship between the viruses hosts in our data set?

So next we download 8 segment 4 viruses hosted by **swine**, but also H1 subtype. We add them to **duck** and **human** data set we used, and now we have a total of 24 sequences with 3 different kinds of **host** by all H1 subtype. Here we list the new sequences we added to the data set.

>gi—4585158—gb—AF091308— /**swine**/4 (HA)/H1N1/USA/1930/// Influenza A virus (A/**swine**/Iowa/15/30 (H1N1)) segment 4 hemagglutinin precursor (HA) mRNA, complete cds

>gi—516372—gb—X57492— /**swine**/4 (HA)/H1N1/USA/1930/// Influenza A virus (A/**swine**/Iowa/15/30(H1N1)) HA1 and HA2 genes for haemagglutinin, genomic RNA

>gi—473536—gb—D00837— /**swine**/4 (HA)/H1N1/United Kingdom/1939/// Influenza A virus (A/**swine**/Cambridge/1939(H1N1)) gene for haemagglutinin, complete cds

>gi—216409465—gb—AB434408— /**swine**/4 (HA)/H1N2/Japan/1980/// Influenza A virus (A/**swine**/Ehime/1/1980(H1N2)) HA gene for hemagglutinin, complete cds

>gi—3831762—gb—AF085413— /**swine**/4 (HA)/H1N2/United Kingdom/1994/// Influenza A virus (A/**swine**/Scotland/410440/94(H1N2)) haemagglutinin precursor, mRNA, complete cds.

>gi—3831764—gb—AF085414— /**swine**/4 (HA)/H1N2/United Kingdom/1994/// Influenza A virus (A/**swine**/England/438207/94(H1N2)) haemagglutinin precursor, mRNA, complete cds.

>gi—3831766—gb—AF085415— /**swine**/4 (HA)/H1N2/United Kingdom/1995/// Influenza A virus (A/**swine**/England/690421/95(H1N2)) haemagglutinin precursor, mRNA, complete cds.

>gi—3831768—gb—AF085416— /**swine**/4 (HA)/H1N2/United Kingdom/1996/// Influenza A virus (A/**swine**/England/17394/96(H1N2)) haemagglutinin precursor, mRNA, complete cds.

And also for ease of presentation, we add the character "p" before the viruses which are hosted by **swine**. As Figure 4.5 and 4.6 show, the viruses hosted by **human**, **duck** and **swine** are clearly and correctly separated into three clusters. As we discussed in the last section, the results are also depend on the data we choose, limited viruses

sequences and only several kinds of hosts are not enough to be used to verify or explain the evolution in biology. The results only can explain the relation between the **hosts** or **subtypes** we used in these data sets. If we want to know what the biological evolution really like, we need much more data and will cost much more.

Finally, to verify our results we also made the high quality program unrooted binary tree by the CompLearn Toolkit, and we get the following clustering results(see cf. Figures 4.7 and 4.8). Even the clustering results are really good, it took hours for the (24×24) distance matrix to output the results.

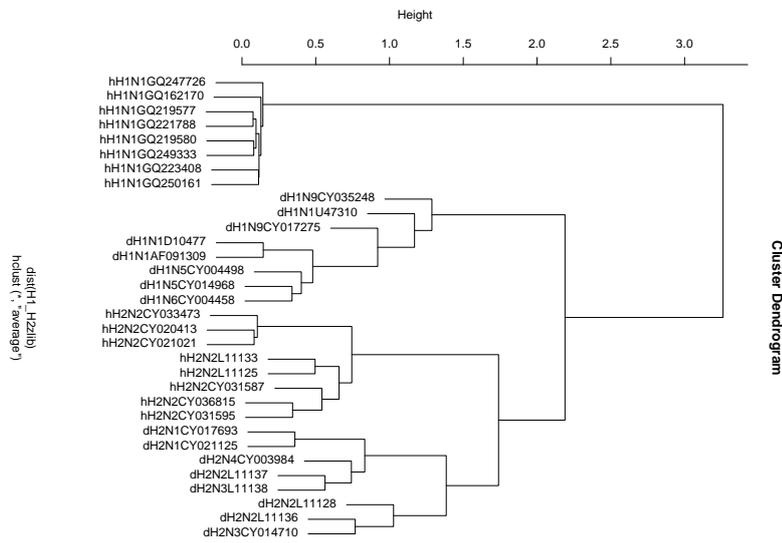


Figure 4.3: Classification of HA sequences hosted by human and duck compr.: zlib

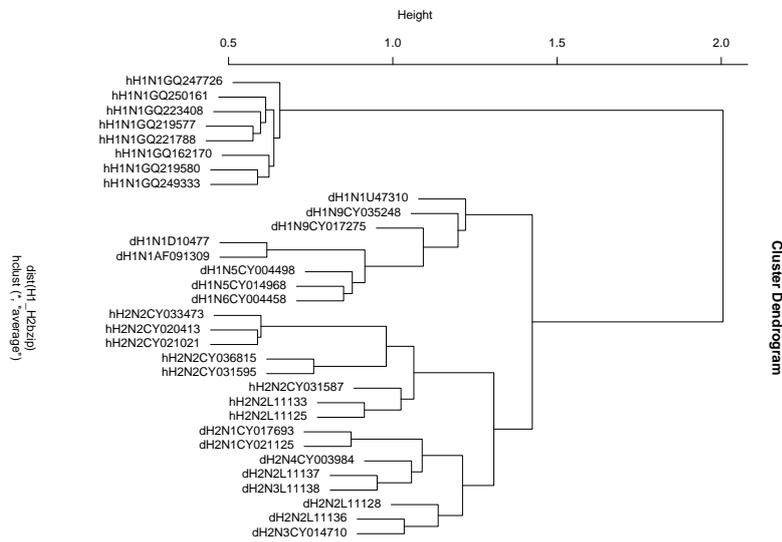


Figure 4.4: Classification of HA sequences hosted by human and duck ; compr.: bzip

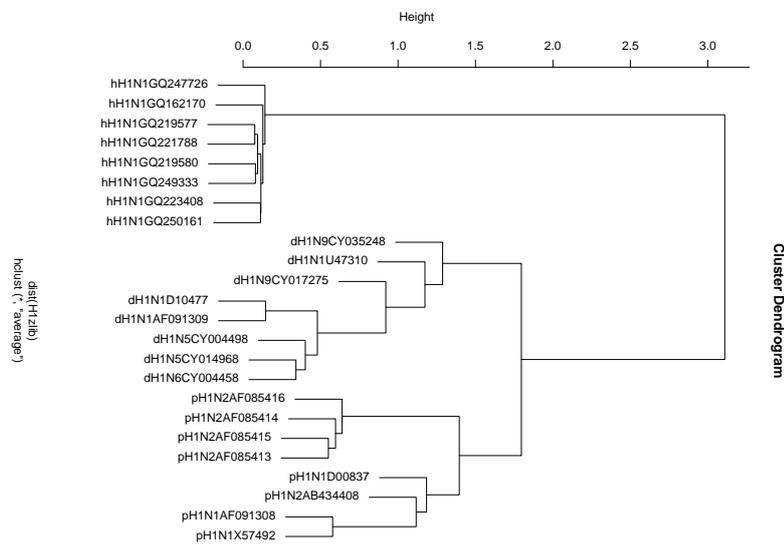


Figure 4.5: Classification of H1 sequences hosted by human, duck and swine; compr.: zlib

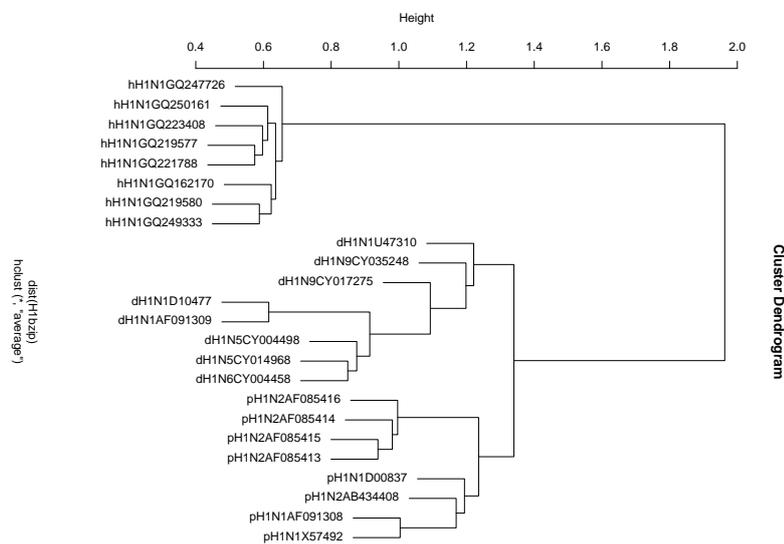


Figure 4.6: Classification of H1 sequences hosted by human, duck and swine; compr.: bzip

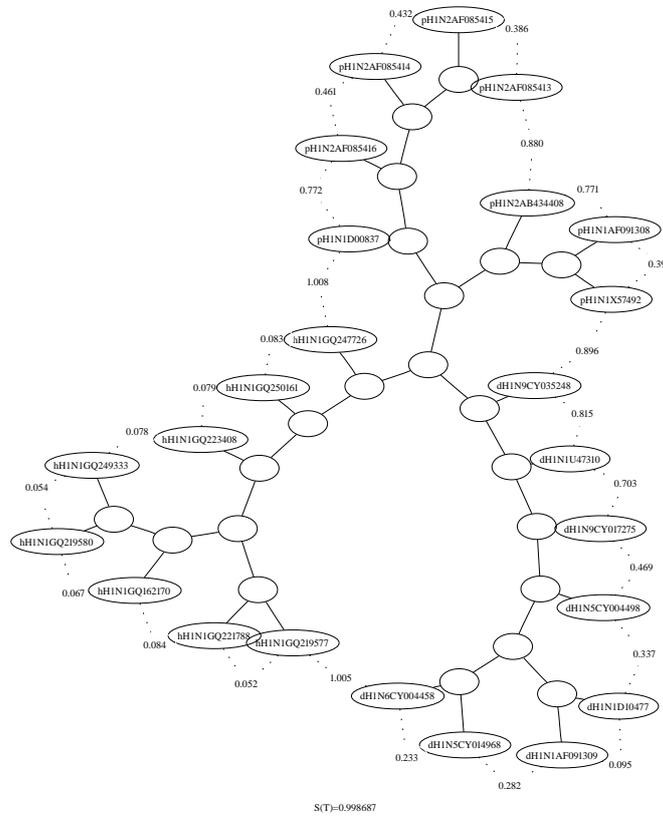


Figure 4.7: Clustering viruses hosted by swine, human and duck via hclust; compr.: zlib

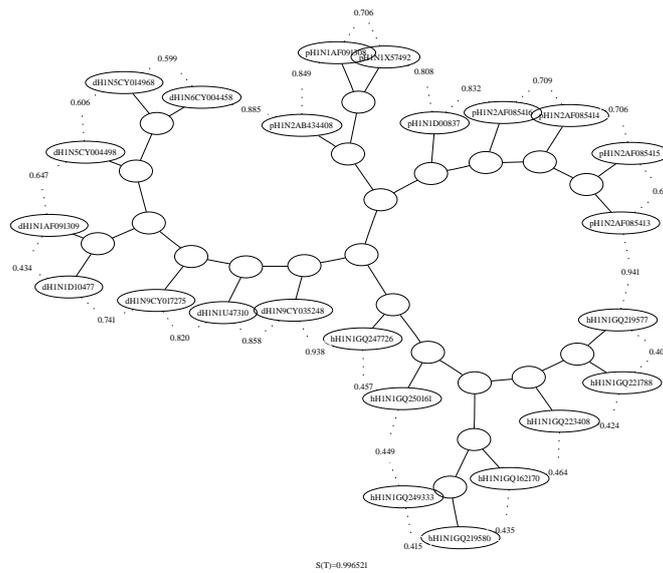


Figure 4.8: Clustering viruses hosted by swine, human and duck via hclust; compr.: bzip

Chapter 5

Conclusions

The usefulness of the normalized compression distance for clustering the HA type of virus data for the HA gene for it (segment 4) has been demonstrated. Though we just used the built-in compressors `zlib` and `bzip` the results are (almost) correct when clustering the resulting distance matrix for the whole data set with `hclust` or spectral clustering via `kLines`. What is also remarkable in this context is the robustness with respect to the completeness of the data. As mentioned above, some data contain only a partial cds but this did not influence the quality of the clustering as the results, e.g., H1N1U47310 and H3N2D21171 have only 1000 letters.

We have not reported the running time here, since it is still in the range of several seconds. Though the quartet tree algorithm by Cilibrasi and Vitányi [9] returns a high quality classification, it lacks scalability, since it tries to optimize a quality function, a task which is NP-hard. So, even for the small example including the 24 data for H1, H2, and H3 resulting in (24×24) distance matrix, it took hours to find the resulting (very good) clustering. In contrast, the clustering algorithms used in this study scale nicely at least up to the amount of data for which the distance matrix is efficiently computable, since they have almost the same running time as the `ncd` algorithm.

Acknowledgement

I would like to express my sincere appreciation to my supervisor, Prof. Thomas Zeugmann, who has significantly improved both my research and this thesis. I also would like to thank Associate Professor Kimihito Ito, who gave precious comments on biology field, for his many helpful suggestions, which has significantly helped my research. Finally, I am indebted to my family for their love and support.

Bibliography

- [1] GNU Octave. <http://www.gnu.org/software/octave/>.
- [2] The R project for statistical computing. <http://www.r-project.org/>.
- [3] Dario Benedetto, Emanuele Caglioti, and Vittorio Loreto. Language trees and zipping. *Phys. Rev. Lett.*, 88(4):048702–1–048702–4, 2002.
- [4] Charles H. Bennett, Péter Gács, Ming Li, Paul M. B. Vitányi, and Wojciech H. Zurek. Information distance. *IEEE Transactions on Information Theory*, 44(4):1407–1423, 1998.
- [5] R. Cilibrasi and P. Vitányi. Automatic meaning discovery using Google. Manuscript, CWI, Amsterdam, 2006.
- [6] R. Cilibrasi and Paul M. B. Vitányi. Clustering by compression. *IEEE Transactions on Information Theory*, 51(4):1523–1545, 2005.
- [7] Rudi Cilibrasi. The compLearn toolkit, 2003-. <http://www.complearn.org/>.
- [8] Rudi Cilibrasi and Paul Vitanyi. Similarity of objects and the meaning of words. In *Theory and Applications of Models of Computation, Third International Conference, TAMC 2006, Beijing, China, May 2006, Proceedings*, volume 3959 of *Lecture Notes in Computer Science*, pages 21–45, Berlin, 2006. Springer.
- [9] Rudi Cilibrasi and Paul M.B. Vitányi. A new quartet tree heuristic for hierarchical clustering. In Dirk V. Arnold, Thomas Jansen, Michael D. Vose, and Jonathan E. Rowe, editors, *Theory of Evolutionary Algorithms*, number 06061 in Dagstuhl Seminar Proceedings. Internationales Begegnungs- und Forschungszentrum für Informatik (IBFI), Schloss Dagstuhl, Germany, 2006.
- [10] Igor Fischer and Jan Poland. New methods for spectral clustering. Technical Report IDSIA-12-04, IDSIA / USI-SUPSI, Manno, Switzerland, 2004.
- [11] National Center for Biotechnology Information. Influenza Virus Resource, information, search and analysis. <http://www.ncbi.nlm.nih.gov/genomes/FLU/FLU.html>.

- [12] Justin Bahl Samantha J. Lycett Michal Worobey Oliver G. Pybus Siu Kit Ma Chung Lam Cheung Jayna Raghvani Samir Bhatt J.S.Malik Peiris Yi Guan Anfrew Rambaut Gavin J. D. Smith, Dhanasekaran Vijaykrishna. Origins and evolutionary genomics of the 2009 swine-origin h1n1 influenza a epidemic. *Nature*, 2009.
- [13] hoehling AA. The great epodemic. *Boston: Little, Brown Co.*, 1961.
- [14] Eamonn Keogh, Stefano Lonardi, and Chotirat Ann Ratanamahatana. Towards parameter-free data mining. In *KDD '04: Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 206–215. ACM Press, 2004.
- [15] Ming Li, Xin Chen, Xin Li, Bin Ma, and Paul M.B. Vitányi. The similarity metric. *IEEE Transactions on Information Theory*, 50(12):3250–3264, 2004.
- [16] Ming Li and Paul Vitányi. *An Introduction to Kolmogorov Complexity and its Applications*. Springer, 3rd edition, 2008.
- [17] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to Information Retrieval*. Cambridge University Press, 2008.
- [18] Peter Palese and Megan L. Shaw. Orthomyxoviridae: The viruses and their replication. In David M. Knipe and Peter M. Howley et al., editors, *Fields' Virology*, pages 1647–1689. Lippincott Williams & Wilkins, Philadelphia, fifth edition, 2007.
- [19] P. Perona and W. Freeman. A factorization approach to grouping. In Hans Burkhardt and Bernd Neumann, editors, *5th European Conference on Computer Vision Freiburg, Germany, June, 2–6, 1998, Proceedings, Volume I*, Lecture Notes in Computer Science, pages 655–670. Springer, 1998.
- [20] Jan Poland and Thomas Zeugmann. Clustering pairwise distances with missing data: Maximum cuts versus normalized cuts. In *Discovery Science, 9th International Conference, DS 2006, Barcelona, Spain, October 2006, Proceedings*, volume 4265 of *Lecture Notes in Artificial Intelligence*, pages 197–208. Springer, 2006.
- [21] Jan Poland and Thomas Zeugmann. Clustering the google distance with eigenvectors and semidefinite programming. In *Knowledge Media Technologies, First International Core-to-Core Workshop*, volume 21 of *Diskussionsbeiträge, Institut für Medien und Kommunikationswissenschaft*, pages 61–69. Technische Universität Ilmenau, 2006.
- [22] Daniel A. Spielman and Shang-Hua Teng. Spectral partitioning works: Planar graphs and finite element meshes. In *Proceedings of the 37th Annual IEEE Conference on Foundations of Computer Science*, pages 96–105. IEEE Computer Society, 1996.

- [23] Paul M. B. Vitányi, Frank J. Balbach, Rudi L. Cilibrasi, and Ming Li. Normalized information distance. In *Information Theory and Statistical Learning*, pages 45–82. Springer, New York, 2008.
- [24] Ulrike von Luxburg. A tutorial on spectral clustering. *Statistics and Computing*, 17(4):395–416, 2007.
- [25] Peter F. Wright, Gabriele Neumann, and Yoshihiro Kawaoka. Orthomyxoviruses. In David M. Knipe and Peter M. Howley et al., editors, *Fields' Virology*, pages 1691–1740. Lippincott Williams & Wilkins, Philadelphia, fifth edition, 2007.
- [26] Stella X. Yu and Jianbo Shi. Multiclass spectral clustering. In *Proceedings of the Ninth IEEE International Conference on Computer Vision*, volume 2, pages 313–319. IEEE Computer Society, 2003.

Appendix

Here we list the description for the viruses dataset(hosted by "duck") we used in our experiments.

>gi—4585160—gb—AF091309— /Avian/4 (HA)/H1N1/Canada/1976/// Influenza A virus (A/duck/Alberta/35/76(H1N1)) segment 4 hemagglutinin precursor (HA) mRNA, complete cds.

>gi—221299—gb—D10477— /Avian/4 (HA)/H1N1/Canada/1976/// Influenza A virus (A/mallard/Alberta/35/1976(H1N1)) gene for haemagglutinin, complete cds.

>gi—1912350—gb—U47310— /Avian/4 (HA)/H1N1/Canada/1976/// Influenza A virus (A/duck/Alberta/35/76(H1N1)) hemagglutinin precursor (HA) mRNA, partial cds.

>gi—78095837—gb—CY004498— /Avian/4 (HA)/H1N5/Canada/1981/08/19/ Influenza A virus (A/pintail duck/ALB/631/1981(H1N5)) segment 4, complete sequence

>gi—115279042—gb—CY014968— /Avian/4 (HA)/H1N5/USA/1978/08/01/ Influenza A virus (A/mallard duck/New York/6861/1978(H1N5)) segment 4, complete sequence

>gi—78095742—gb—CY004458— /Avian/4 (HA)/H1N6/Canada/1977/08/02/ Influenza A virus (A/mallard duck/ALB/42/1977(H1N6)) segment 4, complete sequence

>gi—117572843—gb—CY017275— /Avian/4 (HA)/H1N9/USA/1987/10/19/ Influenza A virus (A/mallard/Ohio/265/1987(H1N9)) segment 4, complete sequence

>gi—209866000—gb—CY035248— /Avian/4 (HA)/H1N9/USA/2007/09/12/ Influenza A virus (A/mallard duck/Minnesota/Sg-00113/2007(H1N9)) hemagglutinin (HA) gene, partial cds

>gi—119365124—gb—CY017693— /Avian/4 (HA)/H2N1/USA/1986/10/20/ Influenza A virus (A/mallard/Ohio/30/1986(H2N1)) segment 4, complete sequence

>gi—134047655—gb—CY021125— /Avian/4 (HA)/H2N1/USA/1986/// Influenza A virus (A/mallard/Ohio/37/1986(H2N1)) segment 4, complete sequence

>gi—408520—gb—L11128— /Avian/4 (HA)/H2N2/Hong Kong/1978/// Influenza A virus (A/duck/Hong Kong/273/78 (H2N2)) hemagglutinin (HA) gene, complete cds.

>gi—408552—gb—L11136— /Avian/4 (HA)/H2N2/USA/1961/// Influenza A virus (A/mallard/MT/Y61 (H2N2)) hemagglutinin (HA) gene, complete cds.

>gi—408598—gb—L11137— /Avian/4 (HA)/H2N2/USA/1978/// Influenza A virus (A/mallard/NY/6750/78 (H2N2)) hemagglutinin (HA) gene, complete cds

>gi—115278357—gb—CY014710— /Avian/4 (HA)/H2N3/Germany/1973/// Influenza A virus (A/duck/Germany/1215/1973(H2N3)) segment 4, complete sequence

>gi—408600—gb—L11138— /Avian/4 (HA)/H2N3/Canada/1976/// Influenza A virus (A/mallard/Ontario/56/76 (H2N3)) hemagglutinin (HA) gene, complete cds.

>gi—78059446—gb—CY003984— /Avian/4 (HA)/H2N4/Canada/2002/08/02/ Influenza A virus (A/mallard/Alberta/149/2002(H2N4)) segment 4, complete sequence

>gi—82652711—gb—CY005943— /Avian/4 (HA)/H3N1/Canada/1976/08/11/ Influenza A virus (A/mallard duck/ALB/26/1976(H3N1)) segment 4, complete sequence

>gi—116235390—gb—AB277754— /Avian/4 (HA)/H3N2/Japan/1977/// Influenza A virus (A/duck/Hokkaido/5/1977(H3N2)) segment 4, complete sequence.

>gi—82654664—gb—CY006026— /Avian/4 (HA)/H3N2/Hong Kong/1975/// Influenza A virus (A/duck/Hong Kong/7/1975(H3N2)) segment 4, complete sequence

>gi—418070—gb—D21171— /Avian/4 (HA)/H3N2/Hong Kong/1977/// Influenza A virus (A/duck/Hong Kong/245/1977(H3N2)) gene for hemagglutinin, partial cds

>gi—189312953—gb—EU74652— /Avian/4 (HA)/H3N2/USA/1978/// Influenza A virus (A/duck/NY/6874/1978(H3N2)) hemagglutinin (HA) gene, complete cds.

>gi—324409—gb—M73771— /Avian/4 (HA)/H3N2/Canada/1976/// Influenza A virus (A/duck/Alberta/78/1976(H3N2)) hemagglutinin gene, complete cds.

>gi—125490300—gb—AB292410— /Avian/4 (HA)/H3N3/Hong Kong/1976/// Influenza A virus (A/duck/Hong Kong/22A/1976(H3N3)) HA gene for haemagglutinin, complete cds.

>gi—82652561—gb—CY005936— /Avian/4 (HA)/H3N3/Canada/1978/08/15/ Influenza A virus (A/mallard duck/ALB/712/1978(H3N3)) segment 4, complete sequence

>gi—125490288—gb—AB292404— /Avian/4 (HA)/H4N1/Hong Kong/1980/// Influenza A virus (A/duck/Hong Kong/951/1980(H4N1)) HA gene for haemagglutinin, complete cds.

>gi—49357153—gb—AY633284— /Avian/4 (HA)/H4N1/Canada/1998/// Influenza A virus (A/mallard/Alberta/47/98(H4N1)) hemagglutinin precursor (HA) gene, partial cds

>gi—82653122—gb—CY005968— /Avian/4 (HA)/H4N1/Canada/1977/08/28/ Influenza A virus (A/mallard duck/ALB/291/1977(H4N1)) segment 4, complete sequence

>gi—82652863—gb—CY005953— /Avian/4 (HA)/H4N2/Canada/1978/08/10/ Influenza A virus (A/mallard duck/ALB/354/1978(H4N2)) segment 4, complete sequence

>gi—82652975—gb—CY005959— /Avian/4 (HA)/H4N2/Canada/1984/08/06/ Influenza A virus (A/mallard duck/ALB/630/1984(H4N2)) segment 4, complete sequence

>gi—82654737—gb—CY006030— /Avian/4 (HA)/H4N2/Hong Kong/1976/// Influenza A virus (A/duck/Hong Kong/24/1976(H4N2)) segment 4, complete sequence

>gi—125490296—gb—AB292408— /Avian/4 (HA)/H4N3/Hong Kong/1977/// Influenza A virus (A/duck/Hong Kong/229/1977(H4N3)) HA gene for haemagglutinin, complete cds.

>gi—195183810—gb—AB450446— /Avian/4 (HA)/H4N3/Mongolia/2007/// Influenza A virus (A/duck/Mongolia/274/2007(H4N3)) HA gene for haemagglutinin, complete cds.

>gi—4240447—gb—AF082040— /Avian/4 (HA)/H5N1/USA/1981/// Influenza A virus (A/duck/Minnesota/1525/81(H5N1)) hemagglutinin H5 mRNA, partial cds.

>gi—1840071—gb—U79453— /Avian/4 (HA)/H5N1/USA/1975/// Influenza A virus (A/mallard/Wisconsin/428/75(H5N1)) hemagglutinin mRNA, partial cds.

>gi—82623227—gb—CY005918— /Avian/4 (HA)/H5N2/Canada/1976/08/12/ Influenza A virus (A/mallard duck/ALB/57/1976(H5N2)) segment 4, complete sequence

>gi—193877643—gb—EU743293— /Avian/4 (HA)/H5N2/USA/1975/// Influenza A virus (A/mallard/WI/42/1975(H5N2)) segment 4 hemagglutinin (HA) gene, complete cds.

>gi—148532723—gb—EF597247— /Avian/4 (HA)/H5N3/Hong Kong/1976/// Influenza A virus (A/duck/Hong Kong/23/1976(H5N3)) hemagglutinin (HA) gene, partial cds

>gi—1840069—gb—U79452— /Avian/4 (HA)/H5N3/USA/1975/// Influenza A virus (A/mallard/Wisconsin/169/75(H5N3)) hemagglutinin mRNA, partial cds

>gi—157169393—gb—EF607888— /Avian/4 (HA)/H5N5/USA/2000/// Influenza A virus (A/mallard/MN/105/2000(H5N5)) hemagglutinin gene, complete cds

>gi—4240449—gb—AF082041— /Avian/4 (HA)/H5N6/Germany/1984/// Influenza A virus (A/duck/Potsdam/2216-4/1984(H5N6)) hemagglutinin H5 mRNA, partial cds.

>gi—87246912—gb—DQ376618— /Avian/4 (HA)/H6N1/Taiwan/1972/// Influenza A virus (A/duck/Taiwan/0526/72(H6N1)) hemagglutinin (HA) gene, complete cds

>gi—193876585—gb—EU743286— /Avian/4 (HA)/H6N1/USA/1969/// Influenza A virus (A/duck/PA/486/1969(H6N1)) segment 4 hemagglutinin (HA) gene, complete cds.

>gi—18074894—gb—AJ410541— /Avian/4 (HA)/H6N2/Hong Kong/1977/// Influenza A virus genomic RNA for haemagglutinin (ha gene) strain A/duck/Hong Kong/134/77 (H6N2)

>gi—115278142—gb—CY014616— /Avian/4 (HA)/H6N2/Hong Kong/1977/// Influenza A virus (A/duck/Hong Kong/d134/1977(H6N2)) segment 4, complete sequence

>gi—78072375—gb—CY004202— /Avian/4 (HA)/H6N3/Canada/1985/08/13/ Influenza A virus (A/mallard duck/ALB/76/1985(H6N3)) segment 4, complete sequence

>gi—78093479—gb—CY004234— /Avian/4 (HA)/H6N3/Canada/1990/08/01/ Influenza A virus (A/mallard duck/ALB/191/1990(H6N3)) segment 4, complete sequence

>gi—78070027—gb—CY004086— /Avian/4 (HA)/H6N4/Canada/1979/08/25/ Influenza A virus (A/pintail duck/ALB/1343/1979(H6N4)) segment 4, complete sequence

>gi—70608906—gb—DQ021680— /Avian/4 (HA)/H6N4/USA/1998/// Influenza A virus (A/mallard/MD/R326/98(H6N4)) hemagglutinin gene, partial cds

>gi—146219366—gb—AB268557— /Avian/4 (HA)/H7N1/Mongolia/2001/// Influenza A virus (A/duck/Mongolia/47/2001(H7N1)) HA gene for haemagglutinin, complete cds

>gi—195926997—gb—AB269694— /Avian/4 (HA)/H7N1/Japan/2003/// Influenza A virus (A/duck/Hokkaido/143/2003(H7N1)) HA gene for haemagglutinin, complete cds

>gi—218436727—gb—AB473543— /Avian/4 (HA)/H7N1/Mongolia/2002/// Influenza A virus (A/duck/Mongolia/867/2002(H7N1)) genomic RNA, segment 4, complete sequence

>gi—115278134—gb—CY014612— /Avian/4 (HA)/H7N1/China/1992/// Influenza A virus (A/duck/Nanchang/1904/1992(H7N1)) segment 4, complete sequence

>gi—138395665—gb—CY021557— /Avian/4 (HA)/H7N1/Italy/2000/01/24/ Influenza A virus (A/duck/Italy/551/2000(H7N1)) segment 4, complete sequence

>gi—62910860—gb—DQ003216— /Avian/4 (HA)/H7N1/Hong Kong/1972/// Influenza A virus (A/duck/Hongkong/301/72(H7N1)) hemagglutinin (HA) gene, complete cds.

>gi—146350799—gb—AB302789— /Avian/4 (HA)/H7N2/Hong Kong/1978/// Influenza A virus (A/duck/Hong Kong/301/1978(H7N2)) HA gene for haemagglutinin, complete cds.

>gi—902756—gb—U20461— /Avian/4 (HA)/H7N2/Hong Kong/1978/// Influenza A virus (A/duck/Hong Kong/293/78(H7N2)) hemagglutinin precursor (HA) mRNA, complete cds.

>gi—115279539—gb—CY015173— /Avian/4 (HA)/H8N2/USA/1991/08// Influenza A virus (A/duck/Alaska/702/1991(H8N2)) segment 4, complete sequence

>gi—195183830—gb—AB450454— /Avian/4 (HA)/H8N4/Japan/1981/// Influenza A virus (A/duck/Hokkaido/95/1981(H8N4)) HA gene for haemagglutinin, complete cds.

>gi—11596270—gb—AF310987— /Avian/4 (HA)/H8N4/Canada/1979/// Influenza A virus (A/Pintail duck/Alberta/114/79(H8N4)) segment 4 hemagglutinin (HA1) mRNA, partial cds.

>gi—11596272—gb—AF310988— /Avian/4 (HA)/H8N4/Canada/1984/// Influenza A virus (A/Mallard duck/Alberta/357/84(H8N4)) segment 4 hemagglutinin (HA1) mRNA, partial cds.

>gi—82653160—gb—CY005970— /Avian/4 (HA)/H8N4/Canada/1977/08/06/ Influenza A virus (A/mallard/Alberta/283/1977(H8N4)) segment 4, complete sequence.

>gi—117163687—gb—EF061122— /Avian/4 (HA)/H8N4/China/2005/// Influenza A virus (A/duck/Yangzhou/02/2005(H8N4)) segment 4, complete sequence

>gi—195183782—gb—AB450435— /Avian/4 (HA)/H8N7/USA/1991/// Influenza A virus (A/duck/Alaska/702/1991(H8N7)) HA gene for haemagglutinin, complete cds.

>gi—49357111—gb—AY633116— /Avian/4 (HA)/H9N1/Canada/1983/// Influenza A virus (A/mallard/Alberta/743/83(H9N1)) hemagglutinin precursor (HA) gene, complete cds

>gi—78096059—gb—CY004642— /Avian/4 (HA)/H9N1/Canada/1983/08/12/ Influenza A virus (A/mallard duck/ALB/506/1983(H9N1)) segment 4, complete sequence

>gi—82623259—gb—CY005919— /Avian/4 (HA)/H9N1/Canada/1983/08/07/ Influenza A virus (A/mallard duck/ALB/396/1983(H9N1)) segment 4, complete sequence

>gi—6007003—gb—AF156382— /Avian/4 (HA)/H9N2/Hong Kong/1977/// Influenza A virus (A/duck/Hong Kong/168/77(H9N2)) segment 4 hemagglutinin precursor, gene, partial cds.

>gi—31339421—gb—AF523386— /Avian/4 (HA)/H9N2/Hong Kong/1976/// Influenza A virus (A/duck/Hong Kong/86/76(H9N2)) hemagglutinin (HA) gene, partial cds

>gi—47157066—gb—AY603067— /Avian/4 (HA)/H9N2/China/// Influenza A virus (A/duck/China(H9N2)) hemagglutinin (HA) gene, complete cds.

>gi—145284497—gb—EF541419— /Avian/4 (HA)/H9N3/Viet Nam/2001/// Influenza A virus (A/duck/Viet Nam/68/2001(H9N3)) segment 4 hemagglutinin (HA) gene, partial cds.

>gi—145284499—gb—EF541420— /Avian/4 (HA)/H9N3/Viet Nam/2001/// Influenza A virus (A/duck/Viet Nam/340/2001(H9N3)) segment 4 hemagglutinin (HA) gene, partial cds.

>gi—113531192—gb—AB271117— /Avian/4 (HA)/H10N1/Hong Kong/1980/// Influenza A virus (A/duck/Hong Kong/938/80(H10N1)) HA gene for haemagglutinin, complete cds.

>gi—82653957—gb—CY005997— /Avian/4 (HA)/H10N1/Canada/1995/08/01/ Influenza A virus (A/mallard/ALB/5/1995(H10N1)) segment 4, complete sequence

>gi—195183803—gb—AB450443— /Avian/4 (HA)/H10N2/Japan/2007/// Influenza A virus (A/duck/Hokkaido/W87/2007(H10N2)) HA gene for haemagglutinin, complete cds.

>gi—125490304—gb—AB292412— /Avian/4 (HA)/H10N3/Hong Kong/1979/// Influenza A virus (A/duck/Hong Kong/786/1979(H10N3)) HA gene for haemagglutinin, complete cds.

>gi—82653900—gb—CY005994— /Avian/4 (HA)/H10N3/Canada/1978/08/08/ Influenza A virus (A/blue-winged teal/ALB/778/1978(H10N3)) segment 4, complete sequence

>gi—114651172—gb—AB274041— /Avian/4 (HA)/H10N4/Japan/2000/// Influenza A virus (A/duck/Hokkaido/18/00(H10N4)) HA gene for haemagglutinin, complete cds.

>gi—195183828—gb—AB450453— /Avian/4 (HA)/H10N5/Japan/2004/// Influenza A virus (A/duck/Hokkaido/24/04(H10N5)) HA gene for haemagglutinin, complete cds.

>gi—195183835—gb—AB450456— /Avian/4 (HA)/H10N5/Mongolia/2003/// Influenza A virus (A/duck/Mongolia/149/03(H10N5)) HA gene for haemagglutinin, complete cds.

>gi—195183823—gb—AB450451— /Avian/4 (HA)/H11N1/Japan/1977/// Influenza A virus (A/duck/Miyagi/47/1977(H11N1)) HA gene for haemagglutinin, complete cds.

>gi—119365309—gb—CY017765— /Avian/4 (HA)/H11N1/USA/1986/11/06/ Influenza A virus (A/black duck/Ohio/194/1986(H11N1)) segment 4, complete sequence

>gi—82654071—gb—CY006003— /Avian/4 (HA)/H11N2/Canada/1991/08/26/ Influenza A virus (A/mallard/ALB/124/1991(H11N2)) segment 4, complete sequence

>gi—82654052—gb—CY006002— /Avian/4 (HA)/H11N3/Canada/1983/08/23/ Influenza A virus (A/mallard duck/ALB/797/1983(H11N3)) segment 4, complete sequence

>gi—119502448—gb—CY018015— /Avian/4 (HA)/H11N3/USA/1986/10/24/ Influenza A virus (A/mallard/Ohio/102/1986(H11N3)) segment 4, complete sequence

>gi—120169268—gb—AB288845— /Avian/4 (HA)/H11N6/United Kingdom/1956/// Influenza A virus (A/duck/England/1/1956(H11N6)) HA gene for haemagglutinin, complete cds.

>gi—115278283—gb—CY014679— /Avian/4 (HA)/H11N6/United Kingdom/1956/// Influenza A virus (A/duck/England/1956(H11N6)) segment 4, complete sequence

>gi—221307—gb—D90306— /Avian/4 (HA)/H11N6/United Kingdom/1956/// Influenza A virus (A/duck/England/1/1956(H11N6)) gene for hemagglutinin precursor, complete cds.

>gi—11596278—gb—AF310991— /Avian/4 (HA)/H12N1/Canada/1983/// Influenza A virus (A/Mallard duck/Alberta/342/83(H12N1)) segment 4 hemagglutinin (HA1) mRNA, partial cds.

>gi—82654128—gb—CY006006— /Avian/4 (HA)/H12N1/Canada/1983/08/06/ Influenza A virus (A/mallard duck/Alberta/342/1983(H12N1)) segment 4, complete sequence.

>gi—112382771—gb—DQ787811— /Avian/4 (HA)/H12N2/Russia/2002/// Influenza A virus (A/duck/Primorie/3691/02(H12N2)) hemagglutinin precursor (HA) gene, partial cds.

>gi—119943225—gb—AB288334— /Avian/4 (HA)/H12N5/Canada/1976/// Influenza A virus (A/duck/Alberta/60/1976(H12N5)) HA gene for haemagglutinin, complete cds.

>gi—120169264—gb—AB288843— /Avian/4 (HA)/H12N5/Japan/2001/// Influenza A virus (A/duck/Hokkaido/66/01(H12N5)) HA gene for haemagglutinin, complete cds.

>gi—134048125—gb—CY021293— /Avian/4 (HA)/H12N5/USA/2005/// Influenza A virus (A/mallard/Maryland/1135/2005(H12N5)) segment 4, complete sequence

>gi—134048144—gb—CY021301— /Avian/4 (HA)/H12N5/USA/2005/// Influenza A virus (A/mallard/Maryland/1153/2005(H12N5)) segment 4, complete sequence

>gi—221309—gb—D90307— /Avian/4 (HA)/H12N5/Canada/1976/// Influenza A virus (A/duck/Alberta/60/1976(H12N5)) gene for hemagglutinin precursor, complete cds.

>gi—118595865—gb—AB284988— /Avian/4 (HA)/H13N6/Russia/1998/// Influenza A virus (A/duck/Siberia/272/1998(H13N6)) HA gene for haemagglutinin, complete cds.

>gi—118722035—gb—AB285094— /Avian/4 (HA)/H13N6/Russia/1998/// Influenza A virus (A/duck/Siberia/272PF/1998(H13N6)) HA gene for haemagglutinin, complete cds.

>gi—324046—gb—M35996— /Avian/4 (HA)/H14/Russia/1982/// Influenza A/Mallard/Gurjev/ hemagglutinin subtype H14 gene.

>gi—324045—gb—M35997— /Avian/4 (HA)/H14/Russia/1982/// Influenza A/Mallard/Gurjev/ hemagglutinin subtype H14 gene.

>gi—120871775—gb—AB289335— /Avian/4 (HA)/H14N5/Russia/1982/// Influenza A virus (A/mallard/Astrakhan/263/1982(H14N5)) HA gene for haemagglutinin, complete cds.

>gi—115278117—gb—CY014604— /Avian/4 (HA)/H14N5/Russia/1982/// Influenza A virus (A/mallard duck/Astrakhan/263/1982(H14N5)) segment 4, complete sequence

>gi—82791465—gb—CY006032— /Avian/4 (HA)/H15N2/Australia/1983/// Influenza A virus (A/Australian shelduck/Western Australia/1756/1983(H15N2)) segment 4, complete sequence.

>gi—126567438—gb—AB295613— /Avian/4 (HA)/H15N8/Australia/1983/// In-
fluenza A virus (A/duck/Australia/341/83(H15N8)) HA gene for haemagglutinin,
complete cds.

>gi—82654235—gb—CY006009— /Avian/4 (HA)/H15N8/Australia/1983/02/23/
Influenza A virus (A/duck/AUS/341/1983(H15N8)) segment 4, complete sequence

>gi—1226068—gb—L43916— /Avian/4 (HA)/H15N8/Australia/1983/// Influenza
A/duck/Australia/341/83 (H15N8) hemagglutinin mRNA, complete cds.

>gi—161878869—gb—EU148600— /Avian/4 (HA)/H16N3/Russia/1983/// Influenza
A virus (A/mallard/Gurjev/785/83(H16N3)) hemagglutinin precursor (HA) gene,
complete cds