# Prediction of Influenza A Virus Infections in Humans using an Artificial Neural Network Learning Approach

Charalambos Chrysostomou[1*], Harris Partaourides[2] and Huseyin Seker[3]

*Abstract*— The Influenza type A virus can be considered as one of the most severe viruses that can infect multiple species with often fatal consequences to the hosts. The Haemagglutinin (HA) gene of the virus has the potential to be a target for antiviral drug development realised through accurate identification of its sub-types and possible the targeted hosts. In this paper, to accurately predict if an Influenza type A virus has the capability to infect human hosts, by using only the HA gene, is therefore developed and tested. The predictive model follows three main steps; (i) decoding the protein sequences into numerical signals using EIIP amino acid scale, (ii) analysing these sequences by using Discrete Fourier Transform (DFT) and extracting DFT-based features, (iii) using a predictive model, based on Artificial Neural Networks and using the features generated by DFT.

In this analysis, from the Influenza Research Database, 30724, 18236 and 8157 HA protein sequences were collected for Human, Avian and Swine respectively. Given this set of the proteins, the proposed method yielded 97.36% ($\pm$ 0.04%), 97.26% ($\pm$ 0.26%), 0.978 ($\pm$ 0.004), 0.963 ($\pm$ 0.005) and 0.945 ($\pm$ 0.005) for the training accuracy validation accuracy, precision, recall and Mathews Correlation Coefficient (MCC) respectively, based on a 10-fold cross-validation. The classification model generated by using one of the largest dataset, if not the largest, yields promising results that could lead to early detection of such species and help develop precautionary measurements for possible human infections.

*Index Terms*— Artificial Neural Network, Amino Acid Indices, Discrete Fourier Transform (DFT), Hemagglutinin (HA) Protein

## I. INTRODUCTION

The Influenza type A virus can be considered one of the most severe virus that can infect both mammals and birds. The genome of the Influenza virus is composed of eight segments that can encode more than 11 proteins [1]. One of the most important proteins is the Haemagglutinin (HA), which is an essential glycoprotein and a principal surface antigen which is responsible for attaching the virions to hosts, deciding the pathogenicity and virulence [1]. Until now, 18 distinct Influenza A HA subtypes have been identified [2], [3].

The Influenza type A virus continually evolves due to the high mutation rate and the constant changes to its genome.

This constant adaptation usually makes any new strain of virus more pathogenic than the previous. Furthermore, these mutations also provide the virus with the ability to cross the species barrier and may also affect the binding pattern of a virus, with catastrophic consequences to the concerned species [4].

In the literature, previous efforts and analysis have been performed to analyse and characterise the phylogenetic diversity, and discover mechanisms that define the severity and distribution of influenza type A virus [5]–[7]. Additionally, as the authors concluded, classification and characterisation of all the sequences with the proposed methods, was difficult [5]–[7], thus a more advanced method is needed. Computational studies exist that tries to characterise and analyse the Influenza type A with promising results [8]. In the proposed method a computational, Artificial Neural Network (ANN) learning based approach is created to predict if a particular virus has the capability to infect humans, by only analysing the HA protein sequence.

The paper is organised as follows: Section II presents the methods and materials developed and used, while Section III presents the results obtained. Finally, concluding remarks are outlined in Section IV.
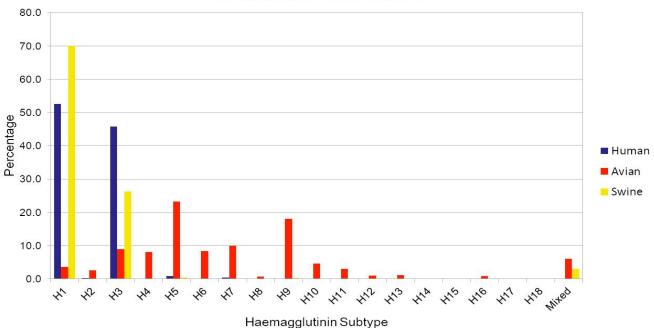
## II. METHODS AND MATERIAL

### A. Influenza A Hemagglutinin Proteins Data Set

For the proposed analysis 57117 HA Influenza type A protein sequences are collected from the Influenza Research Database [9], for three species, Human, Avian and Swine. More specifically, as Table I shows 30724, 18236 and 8157 HA protein sequences were collected for Human, Avian and Swine respectively. Furthermore, Table I shows the specific number of sequences for each class of HA 1-18. Finally, figure 1 illustrates the percentage of HA proteins per class and per species. For the analysis, classification of the HA protein sequences, based on the ability of the virus to infect human hosts, the data were separated into two groups. The first group contained all the sequences from HA 1-18 for the viruses that have the ability to infect the Humans hosts, and the second group with the sequences that have the potential to infect the Avian and Swine hosts. For the first and second groups, the total number of 30724 and 26393 HA protein sequences were used respectively.

### B. Data conversion and Normalisation

In this paper digital signal processing techniques are used to extract information that can be directly used to characterise the HA proteins. In the literature, various methods

[1]Computation-based Science and Technology Research Center, The Cyprus Institute, 20 Konstantinou Kavafi Street, 2121, Aglantzia, Nicosia, Cyprus

[2]Department of Electrical Engineering, Computer Engineering and Informatics, Cyprus University of Technology, 30 Archbishop Kyprianou Str., 3036 Limassol, Cyprus

[3]Department of Computer and Information Sciences, Faculty of Engineering and Environment, The University of Northumbria at Newcastle, NE1 8ST, Newcastle-upon-Tyne, The United Kingdom

c.chrysostomou@cyi.ac.cy, c.partaourides@cut.ac.cy, huseyin.seker@northumbria.ac.uk

*Corresponding Author

Fig. 1.   HA Protein Sequences

TABLE I
NUMBER OF HA PROTEIN SEQUENCES

| HA Subtype | Human | Avian | Swine |
|---|---|---|---|
| H1 | 16145 | 650 | 5714 |
| H2 | 96 | 462 | 2 |
| H3 | 14055 | 1621 | 2138 |
| H4 | 0 | 1478 | 5 |
| H5 | 269 | 4242 | 32 |
| H6 | 0 | 1529 | 2 |
| H7 | 104 | 1809 | 1 |
| H8 | 0 | 121 | 0 |
| H9 | 13 | 3281 | 20 |
| H10 | 4 | 834 | 1 |
| H11 | 0 | 547 | 0 |
| H12 | 0 | 178 | 0 |
| H13 | 0 | 200 | 0 |
| H14 | 0 | 17 | 0 |
| H15 | 0 | 13 | 0 |
| H16 | 0 | 150 | 0 |
| H17 | 0 | 0 | 0 |
| H18 | 0 | 0 | 0 |
| Mixed | 38 | 1104 | 242 |
| Total | 30724 | 18236 | 8157 |

used signal processing in bioinformatics for analysing and characterising protein sequences [10]–[14] such as Complex resonant recognition model in analysing influenza a virus subtype protein sequences [10], CISAPS: Complex informational spectrum for the analysis of protein sequences [13] and Structural classification of protein sequences based on signal processing and support vector machines [14]. Furthermore, previous studies [15] where signal processing was used to analyse influenza A HA proteins aimed to identify new therapeutic targets for drug development by better understanding the interaction of the influenza virus and its receptors.

For the proposed analysis, signal processing methods are used, and more specifically Discrete Fourier Transform (DFT), as shown in equations 1-3. The analysis was performed directly to absolute spectrum. Before applying DFT to the HA protein sequences, Electron-ion interaction potential (EIIP) [16], [17] amino acid index, was used to convert alphanumerical sequences. The complete list of the EIIP amino acid index can be found in Table II.

Discrete Fourier Transform (DFT)

$$X(n) = \sum_{m=0}^{N-1} x(m)e^{-j(2\pi/N)nm} \quad n = 0, 1, ..., N-1 \quad (1)$$

where $X(n)$ are the DFT coefficients, N is the total number of points in the series and $x(m)$ is the $m$th member of the numerical series. As the DFT coefficients contain two mirror parts, only the $(N/2)$ points of the series will be used.
The output of DFT is a complex sequence and can be formulated as

$$X(n) = (R(n) + jI(n)), \quad n = 0, 1, ..., (N-1)/2 \quad (2)$$

where $R(n)$ and $I(n)$ are the Real and Imaginary parts of the sequence, respectively. The absolute spectrum $(S_{(n)})$ can be formulated as

$$S_{(n)} = X(n)X^*(n) = |X(n)|^2, \quad n = 0, 1, ..., (N-1)/2 \tag{3}$$

where $X(n)$ are the DFT coefficients of the series $x(n)$, $X*(n)$ are the complex conjugates.

The coefficients from the absolute spectrum will be used as a feature set to represent the characteristics of different classes of proteins' secondary structure. The HA influenza A virus proteins sequences have different lengths, and zero-padding was used to extend all the protein sequences to $N = 1024$ before applying DFT. After DFT is applied the output of the absolute spectrum includes 513 features. These features are used as input to the ANN model.

TABLE II

EIIP VALUES

| Amino acid | EIIP Values |
|---|---|
| Leucine | 0.0000 |
| Isoleucine | 0.0000 |
| Asparagine | 0.0036 |
| Glycine | 0.0050 |
| Glutamic acid | 0.0057 |
| Valine | 0.0058 |
| Proline | 0.0198 |
| Histidine | 0.0242 |
| Lysine | 0.0371 |
| Alanine | 0.0373 |
| Tyrosine | 0.0516 |
| Tryptophan | 0.0548 |
| Glutamine | 0.0761 |
| Methionine | 0.0823 |
| Serine | 0.0829 |
| Cysteine | 0.0829 |
| Threonine | 0.0941 |
| Phenylalanine | 0.0946 |
| Arginine | 0.0959 |
| Aspartic acid | 0.1263 |

*C. Artificial Neural Network - Experimental Evaluation*

Artificial Neural Networks (ANN) [18] are a computational method, based on an extensive collection of artificial neurons, which mirrors the process a living brain solves problems. Each neuron connects to multiple other neurons, which can enforce or repress the impact on the activation event of the connected neurons. The ANNs are considered as self-trained, rather than explicitly programmed, and employed in research fields where the discovery of features and classification is challenging in traditional classification systems.

For the proposed work, the ANN receives an input of 513 features derived from the influenza type A virus pre-processing and returns the probability of the virus infecting humans. For the proposed work, the binary classification consist of 57117 samples of which 30724 can infect humans.

The network setup consists of a single hidden layer of 128 units, Glorot-style uniform for initialization and rectified linear units for the activation function. In order to train the ANN, the Adam optimizer [19] was used with mini batch size of 128 for 200 epochs. We use 10 fold cross-validation and show the network performance based on accuracy. The model was implemented by utilising the Tensorflow [20] and Keras [21] libraries. A visual representation of the model can be seen in Figure 2.
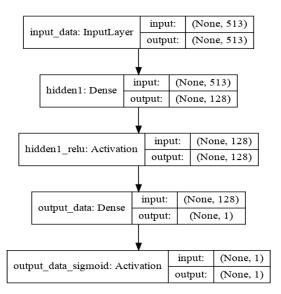


Fig. 2. Artificial Neural Network Model

The high performance of both training and testing shows that for this type of problem more advanced Neural Network models (such as Deep Neural Networks) and regularization techniques are not needed.

## III. RESULTS AND DISCUSSION

In this paper, a classification model is presented, based on Artificial Neural Networks, for the analysis and classification Influenza type A based upon the ability to infect a human host solely by using the HA protein sequence. To ensure that the proposed classification model is accurate and the results can be generalised, 10-fold cross-validation was used. The total accuracy of the predictive model with average training accuracy, testing accuracy, precision, recall and MCC of 97.36% ($\pm$ 0.04%), 97.26% ($\pm$ 0.26%), 0.978 ($\pm$ 0.004), 0.963 ($\pm$ 0.005) and 0.945 ($\pm$ 0.005) for the training accuracy validation accuracy, precision, recall and Mathews Correlation Coefficient (MCC) respectively. As the results show, the proposed model can distinguish HA protein sequences with extremely high accuracy whenever the virus under investigation will have the capability to infect human hosts. Detailed results can be found in Table III.

## IV. CONCLUSIONS

The paper presents a highly successful predictive model to identify and differentiate Influenza type A virus, which can and cannot infect Humans, based on the HA gene, which is considered as a highly potential antiviral drug candidate. The classification model was created by utilising one of the largest possible datasets, if not the largest, in order to better generalise the model. The classification model obtained over

TABLE III

ACCURACY RESULTS FOR THE PREDICTION OF INFLUENZA A VIRUS INFECTIONS

| Fold | Training Accuracy | Validation Accuracy | Validation Precision | Validation Recall | Validation MCC |
|---|---|---|---|---|---|
| 1 | 0.974 | 0.970 | 0.976 | 0.959 | 0.940 |
| 2 | 0.973 | 0.971 | 0.982 | 0.955 | 0.941 |
| 3 | 0.973 | 0.970 | 0.973 | 0.962 | 0.940 |
| 4 | 0.974 | 0.971 | 0.974 | 0.964 | 0.942 |
| 5 | 0.974 | 0.970 | 0.982 | 0.954 | 0.941 |
| 6 | 0.973 | 0.975 | 0.981 | 0.966 | 0.950 |
| 7 | 0.974 | 0.972 | 0.973 | 0.967 | 0.944 |
| 8 | 0.973 | 0.978 | 0.982 | 0.971 | 0.957 |
| 9 | 0.974 | 0.975 | 0.980 | 0.966 | 0.950 |
| 10 | 0.974 | 0.972 | 0.976 | 0.963 | 0.944 |
| Average | 97.36% ($\pm$ 0.04%) | 97.26% ($\pm$ 0.26%) | 0.978 ($\pm$ 0.004) | 0.963 ($\pm$ 0.005) | 0.945 ($\pm$ 0.005) |

the 10-fold cross validation yielded the average training accuracy, testing accuracy, precision, recall and MCC of 97.36% ($\pm$ 0.04%), 97.26% ($\pm$ 0.26%), 0.978 ($\pm$ 0.004), 0.963 ($\pm$ 0.005) and 0.945 ($\pm$ 0.005) respectively. Also, by applying signal processing technique, namely Discrete Fourier Transform, on the protein sequences, it was found that useful spectral characteristic features can be distinguished that are capable of representing the protein groups and thus further enhanced using the Artificial Neural Network based classifier.

As reported in the literature, there are over 600 amino acid indices, where each represents a unique physicochemical feature of the protein [13], in contrast to the one amino acid index used throughout this study. Future studies are required to identify any potential amino acid indices that are capable of representing, characterising and classify the Influenza type A HA protein sequences. A computational tool that is capable of classifying and distinguishing potentially dangerous viruses that have the capability to infect Human hosts will be critical and required to monitor future outbreaks. Future studies will investigate additional machine learning approaches, to explore the efficiency of the methodology, using signal processing methods to encode protein sequences.

## REFERENCES

[1] R. G. Webster, W. J. Bean, O. T. Gorman, T. M. Chambers, and Y. Kawaoka, "Evolution and ecology of influenza a viruses." *Microbiological reviews*, vol. 56, no. 1, pp. 152–179, 1992.

[2] R. A. Fouchier, V. Munster, A. Wallensten, T. M. Bestebroer, S. Herfst, D. Smith, G. F. Rimmelzwaan, B. Olsen, and A. D. Osterhaus, "Characterization of a novel influenza a virus hemagglutinin subtype (h16) obtained from black-headed gulls," *Journal of virology*, vol. 79, no. 5, pp. 2814–2822, 2005.

[3] Y. Wu, Y. Wu, B. Tefsen, Y. Shi, and G. F. Gao, "Bat-derived influenza-like viruses h17n10 and h18n11," *Trends in microbiology*, vol. 22, no. 4, pp. 183–191, 2014.

[4] N. Nunthaboot, T. Rungrotmongkol, M. Malaisree, N. Kaiyawet, P. Decha, P. Sompornpisut, Y. Poovorawan, and S. Hannongbua, "Evolution of human receptor binding affinity of h1n1 hemagglutinins from 1918 to 2009 pandemic influenza a virus," *Journal of chemical information and modeling*, vol. 50, no. 8, pp. 1410–1417, 2010.

[5] J.-M. Chen, Y.-X. Sun, J.-W. Chen, S. Liu, J.-M. Yu, C.-J. Shen, X.-D. Sun, and D. Peng, "Panorama phylogenetic diversity and distribution of type a influenza viruses based on their six internal gene sequences," *Virology Journal*, vol. 6, no. 1, p. 137, 2009.

[6] S. Liu, K. Ji, J. Chen, D. Tai, W. Jiang, G. Hou, J. Chen, J. Li, and B. Huang, "Panorama phylogenetic diversity and distribution of type a influenza virus," *PLoS One*, vol. 4, no. 3, p. e5022, 2009.

[7] W. Shi, F. Lei, C. Zhu, F. Sievers, and D. G. Higgins, "A complete analysis of ha and na genes of influenza a viruses," *PloS one*, vol. 5, no. 12, p. e14454, 2010.

[8] Z. Rehman, R. Zafar, U. Amir, U. H. Niazi, and A. Fahim, "Characterization of evolutionary changes in hemagglutinin of influenza h1n1 virus: a computational analysis," *VirusDisease*, vol. 27, no. 1, pp. 34–40, 2016.

[9] R. B. Squires, J. Noronha, V. Hunt, A. García-Sastre, C. Macken, N. Baumgarth, D. Suarez, B. E. Pickett, Y. Zhang, C. N. Larsen *et al.*, "Influenza research database: an integrated bioinformatics resource for influenza research and surveillance," *Influenza and other respiratory viruses*, vol. 6, no. 6, pp. 404–416, 2012.

[10] C. Chrysostomou, H. Seker, N. Aydin, and P. I. Haris, "Complex resonant recognition model in analysing influenza a virus subtype protein sequences," in *Information Technology and Applications in Biomedicine (ITAB), 2010 10th IEEE International Conference on*. IEEE, 2010, pp. 1–4.

[11] C. J. Carmona, C. Chrysostomou, H. Seker, and M. del Jesus, "Fuzzy rules for describing subgroups from influenza a virus using a multi-objective evolutionary algorithm," *Applied Soft Computing*, vol. 13, no. 8, pp. 3439–3448, 2013.

[12] C. Chrysostomou, H. Seker, and N. Aydin, "Effects of windowing and zero-padding on complex resonant recognition model for protein sequence analysis," in *Engineering in Medicine and Biology Society, EMBC, 2011 Annual International Conference of the IEEE*. IEEE, 2011, pp. 4955–4958.

[13] C. Chrysostomou, H. Seker, and N. Aydin, "Cisaps Complex informational spectrum for the analysis of protein sequences," *Advances in bioinformatics*, vol. 2015, 2015.

[14] C. Chrysostomou and H. Seker, "Structural classification of protein sequences based on signal processing and support vector machines," in *Engineering in Medicine and Biology Society (EMBC), 2016 IEEE 38th Annual International Conference of the*. IEEE, 2016, pp. 3088–3091.

[15] V. Veljkovic, N. Veljkovic, C. Muller, S. Muller, S. Glisic, V. Perovic, and H. Kohler, "Characterization of conserved properties of hemagglutinin of h5n1 and human influenza viruses: possible consequences for therapy and infection control," *BMC Structural Biology*, vol. 9, no. 1, p. 21, 2009.

[16] V. Veljkovic, I. Cosic, B. Dimitrijevic, and D. LalovicC, "Is it possible to analyze DNA and protein sequences by the methods of digital signal processing?" *IEEE Transaction on Biomedical Engineering*, vol. 32, no. 5, pp. 337–341, 1985.

[17] K. Gopalakrishnan, R. Zadeh, K. Najarian, and A. Darvish, "Computational analysis and classification of p53 mutants according to primary structure," in *2004 IEEE Computational Systems Bioinformatics Conference, Proceedings*, 2004, Proceedings Paper, pp. 694–695.

[18] N. Gupta, "Artificial neural network," *Network and Complex Systems*, vol. 3, no. 1, pp. 24–28, 2013.

[19] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *CoRR*, vol. abs/1412.6980, 2014. [Online]. Available: http://arxiv.org/abs/1412.6980

[20] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin *et al.*, "Tensorflow: Large-scale machine learning on heterogeneous distributed systems," *arXiv preprint arXiv:1603.04467*, 2016.

[21] F. Chollet, "Keras," 2015.