

Journal of Mechanics in Medicine and Biology  
Vol. 12, No. 2 (2012) 1240009 (11 pages)  
© World Scientific Publishing Company  
DOI: 10.1142/S021951941240009X



1  
2  
3  
4  
5  
6 **INFLUENZA A SUBTYPING AND HOST ORIGIN**  
7 **CLASSIFICATION USING PROFILE HIDDEN**  
8 **MARKOV MODELS**  
9

10  
11 **FAYROZ F. SHERIF**

12 *Bioelectronics Department*  
13 *Modern University for Technology and Information*  
14 *Cairo, Egypt*  
15 *ffs@k-space.org*

16 **MAHMOUD EL-HEFNAWI**

17 *Informatics and Systems Department*  
18 *National Research Centre, Giza, Egypt*  
19 *mahef@aucegypt.edu*

20 **YASSER M. KADAH**

21 *Biomedical Engineering Department*  
22 *Cairo University, Giza, Egypt*

23  
24 Received

25 Revised

26 Accepted

27  
28 Influenza is one of the most important emerging and reemerging infectious diseases, causing high  
29 morbidity and mortality in communities (epidemic) and worldwide (pandemic). Here, classifica-  
30 tion of human vs. non-human influenza, and subtyping of human influenza viral strains virus  
31 is done based on profile hidden Markov models (HMM). The classical ways of determining  
32 influenza viral subtypes depend mainly on antigenic assays, which is time-consuming and not  
33 fully accurate. The introduced technique is much cheaper and faster, yet usually can still yield  
34 high accuracy. Multiple sequence alignments were done for the 16 HA subtypes and 9 NA  
35 subtypes, followed by profile-HMMs models generation, calibration and evaluation using the  
36 HMMER suite for each group. Subtyping accuracy of all HA and NA models achieved 100%,  
37 while host classification achieved accuracies around 53% and 95.1% according to HA subtype.

38 *Keywords:* Bioinformatics; influenza virus; profile hidden Markov model.

39 **1. Introduction**

40 Influenza A viruses belong to the Orthomyxoviridae family of negative sense, single-  
41 stranded, segmented RNA viruses. The RNA core consists of 8 gene segments.  
42 Immunologically, the most significant surface proteins include Hemagglutinin  
43 HA (16 subtypes) and Neuraminidase NA (9 subtypes). Influenza A subtypes are

*F. F. Sherif, M. El-Hefnawi & Y. M. Kadah*

1 traditionally identified by their HA and NA proteins.<sup>1,2</sup> The HA and NA proteins are  
2 integral membrane proteins and are considered as the major surface antigen of the  
3 influenza virus virion. HA is responsible for binding of virions to host cell receptors  
4 and for fusion between the virion envelope and the host cell.<sup>3</sup> The role of NA is to free  
5 virus particles from host cell receptors, to permit progeny virions to escape from the  
6 cell in which they arose, and so facilitate virus spread.<sup>4</sup> All the 16 subtypes of HA and  
7 9 subtypes of NA are found in avian but the first three subtypes H1, H2, H3 and  
8 recently H5, are found in human influenza viruses.<sup>5</sup> The most common strains which  
9 infect humans during the annual influenza season are H1N1 and H3N2.<sup>6</sup> Swine  
10 influenza is known to be caused by influenza A subtypes H1N1, H1N2, H3N1, and  
11 H3N2. Rapid virus subtype identification is critical for accurate diagnosis of human  
12 infections, effective response to epidemic outbreaks and global-scale surveillance of  
13 highly pathogenic subtypes such as avian influenza H5N1 and H1N1 2009 virus.<sup>7</sup> The  
14 classical ways of subtyping influenza A virus for HA segments are hemagglutination  
15 inhibition (HI) assay which are capable of distinguishing antigenic differences  
16 between influenza even of the same subtype. However, as noted in Ref. 8, when  
17 working with uncharacterized viruses or antibody subtypes, the library of reference  
18 reagents required for identifying antigenically distinct influenza viruses and/or  
19 antibody specificities from multiple lineages of a single HA subtype requires extensive  
20 laboratory support for the production and optimization of reagents. Another possible  
21 method is the subtyping of HA genes by reverse transcription PCR.<sup>9</sup> Real-time PCR  
22 is highly specific. But there are some things to be considered such as cost and time.  
23 While the cost of primers is probably manageable, probes are very expensive. There  
24 will be a lag time as we will have to obtain all the probes and primers and do  
25 validation studies. A common way to find which subtype a genetic sequence belongs  
26 to is through the BLAST search.<sup>10</sup> However, there are issues associated with the  
27 BLAST algorithm as described in Ref. 11. Most importantly, the BLAST result  
28 cannot reveal important mutations that may be functionally related to the structure  
29 and function of proteins.

30 Profile hidden Markov models (HMMs) are statistical models of multiple sequence  
31 alignments.<sup>12</sup> They capture position-specific information about how conserved each  
32 column of the alignment is, and which residues are likely. Recently related studies  
33 have been conducted to classify influenza virus antigenic types and hosts. An Inte-  
34 grated approach of using decision trees and HMM for subtype prediction of human  
35 influenza A virus — HA subtypes (H1, H2 and H3) and NA subtypes (N1 and N2)  
36 — has been introduced in Ref. 13. They extracted some informative positions from  
37 decision tree algorithms in the Weka package, and then modeled into profiles  
38 through hidden Markov modeling at nucleotide level, using HMMER with subtype  
39 prediction accuracy of 88% for human subtypes. Also, they developed a web system  
40 for accurate subtype detection of human influenza virus sequences only. The pre-  
41 liminary experiment showed that this system is easy-to-use but not powerful in  
42 identifying human influenza subtypes and there is no facility to use protein  
43 sequences. Another study in Ref. 14 applied two machine learning techniques

1 (decision trees and support vector machines) to identify the origin of latest pandemic  
2 outbreak of H1N1 viral strains. Their results have shown that human and swine  
3 groups are well distinguishable, with classification accuracy above 95% at prediction.  
4 All sequences from HA, M, NA, NP, NS, PA, PB1, and PB2 are classified as swine  
5 influenza, which means sequences in these segments are more closely related to  
6 Swine strain. Therefore, it was suggested that the latest pandemic viral strain is of  
7 swine origin. Finally, the most recently study discussed in Ref. 15 has applied the  
8 feed-forward back-propagation neural network for the classification analysis of  
9 influenza virus.

10 Our study aims to generalize and extend influenza subtype and host classification  
11 to include all influenza A viral subtypes and host origins, by developing a prediction  
12 tool using Profile HMM at protein level, for identifying all influenza viral strains in  
13 the different hosts not only human. In this work, the subtype prediction achieved  
14 100% accuracy while host origin identification achieved accuracies around 53% and  
15 95.1% according to HA subtype.

## 16 17 18 **2. Data and Methods**

### 19 **2.1. Data collection**

20 All sequences were downloaded from the NCBI's (National Center for Biotechnology  
21 Information) Influenza Virus Resources.<sup>16</sup> We ensured the downloaded sequences  
22 were non-redundant and the complete isolation of HA and NA segments. Part of the  
23 data is used for training and the remaining part is used for testing (Table 1). We used  
24 amino acid sequences because they are known to give more reliable results than  
25 nucleotide sequences when the sequence divergence is high.<sup>17</sup>

### 26 27 28 **2.2. Multiple sequence alignment (MSA)**

29 One of the cornerstones of modern bioinformatics is the comparison or alignment of  
30 protein sequences. Sequences can be aligned across their entire length (global  
31 alignment) or only in certain regions (local alignment).<sup>18</sup> Each group of training sets  
32 found in Table 1 was collectively aligned using Clustal X program, which supports  
33 multiple sequence alignment for protein sequences through window graphical user  
34 interface and built by adding the sequences sequentially to the growing MSA pro-  
35 duced a consensus sequence representing the highly conserved regions from the  
36 aligned sequences.<sup>19,20</sup>

### 37 38 39 **2.3. Modeling using profile HMM**

40 Profile HMM techniques are among the most powerful methods for protein homology  
41 detection scoring them above the noise level.<sup>21</sup> HMM profile includes more flexible  
42 information on a given set of sequences than a single sequence.<sup>22</sup> Therefore, database  
43 search methods using profiles is more sensitive to remote similarities than those

*F. F. Sherif, M. El-Hefnawi & Y. M. Kadah*

Table 1. Number of downloaded sequences used for each subtype of HA and NA segments.

HA segment	Group	# of training sequences	# of test sequences
H1	Total H1	971	100
	H1-Human	749	100
	H1-Avian	105	10
H2	H1-Swine	300	68
	Total H2	126	50
	H2-Human	50	13
H3	H2-Avian	124	30
	Total H3	814	100
	H3-Human	550	69
H4	H3-Avian	263	30
	H3-Swine	100	29
	Total (Avian)	200	64
H5	Total H5	1500	256
	H5-Human	110	33
	H5-Avian	1200	184
H6	Total (Avian)	150	40
H7	Total (Avian)	200	64
H8	Total (Avian)	15	4
H9	Total H9	400	97
	H9-Avian	400	42
	H9-Swine	13	2
H10	Total (Avian)	40	7
H11	Total (Avian)	40	11
H12	Total (Avian)	15	4
H13	Total (Avian)	25	5
H14	Total (Avian)	10	2
H15	Total (Avian)	10	2
H16	Total (Avian)	12	4
NA segment			
N1	Total N1	1500	205
	N1-Human	600	56
	N1-Avian	830	70
N2	N1-Swine	100	20
	Total N2	1500	561
	N2-Human	761	100
N3	N2-Avian	700	124
	N2-Swine	191	40
	Total N3	102	40
N4	Total N4	40	10
N5	Total N5	65	20
N6	Total N1	261	15
N7	Total N7	100	20
N8	Total N8	300	30
N9	Total N9	80	20

based on pairwise alignments (e.g., regular BLAST). In particular, profile HMM have generated good results, and are today employed by several databases such as Pfam and Superfamily.<sup>23,24</sup> We divided our analysis into two main steps; profile HMM model building and database searching.

1 Model building involves converting a multiple alignment of each group of  
2 sequences into a probabilistic model, while database searching involves scoring a  
3 sequence to the profile HMM. One of the most widely used profile HMM packages is  
4 HMMER packages.

#### 5 6 **2.4. Model building**

7 A profile HMM is a probabilistic model of multiple alignments of related proteins. The  
8 alignment is modeled using a series of nodes (roughly one per alignment column) each  
9 composed of three states: match, insert and delete. Match and insert states emit amino  
10 acids with probabilities learned during model estimation while delete states are quiet.  
11 Insertions and deletions with respect to the HMM are modeled by insert and delete  
12 states and transition probabilities to them.<sup>12</sup> “Hmmbuild” program in HMMER  
13 package v2.3.2 was used to build a different HMM profiles for each subtype of HA and  
14 NA segments; the input to “Hmmbuild” program were the pre-aligned sequences of  
15 each group in Table 1. In order to increase the sensitivity of database search we used  
16 “hmmcalibrate” program in HMMER to calculate the E-value. The E-value is quite  
17 literally the expected number of false positives at this raw score; the larger the database  
18 you search, the greater the number of expected false positives. HMM database has been  
19 built by concatenating HMM files that are already built and calibrated.<sup>25</sup>

#### 20 21 **2.5. Database searching**

22 Any sequence can be compared to a model by calculating the probability that the  
23 sequence was generated by that model. The negative logarithm of this probability  
24 corresponds to the NULL score calculated for a simple HMM. To score a match to  
25 HMM we have two algorithms: Viterbi algorithm to give the probability of the most  
26 probable alignment with the sequence or Forward algorithm to give the full prob-  
27 ability of a sequence aligning to the profile HMM.<sup>26</sup> “Hmsearch” program in  
28 HMMER package searches one or more sequences against HMM profile. The output  
29 of the program is the sequence family classification top hits list, ranked by E-value.  
30 The scores and E-values here reflect the confidence that this query sequence contains  
31 one or more domains belonging to a domain family. “Hmmpfam” program Searches  
32 an HMM database for matches to a query sequence and get score for each model.<sup>23</sup>

### 33 34 35 **3. Results**

36 Multiple sequence alignments were done for the 16 HA subtypes, 9 NA subtypes and 12  
37 “HA-Host” host specific subtypes, using ClustalX, followed by profile-HMMs models  
38 building, calibration and database generation using the HMMER suite for each group.

#### 39 40 41 **3.1. Subtyping classification results**

42 Subtyping classification was done by scoring the entire test-sets (human) (Table 1),  
43 with each HA and NA HMM models, using “Hmmpfam” program in HMMER suite.

*F. F. Sherif, M. El-Hefnawi & Y. M. Kadah*

Table 2. Summary of host classification results of influenza A virus using HMMs.

HA subtype	Host	Accuracy	Sensitivity	Specificity
H1	Human	94.4%	93.7%	95.7%
	Avian	89.5%	100%	95.3%
	Swine	84.5%	90.3%	82.9%
H2	Human	95.1%	100%	96%
	Avian	90%	91.7%	87.5%
H3	Human	80.8%	86.9%	71.1%
	Avian	90.9%	82.4%	92.7%
	Swine	78.7%	71.4%	78.8%
H5	Human	53%	95.2%	43.5%
	Avian	63%	58.1%	76.9%
H9	Avian	55%	46.7%	80%
	Swine	90%	80%	93.3%

Matches to the right HA or NA subtype were classified as true hits. Matches to a different subtype were classified as false hits. The accuracies of classification results achieved 100%. These results are encouraging and bear great promise for application to influenza virus classification. Therefore any viral strain like H1N1, H1N2, H2N2, H3N2, H5N1 and H9N2 can be accurately classified using HMMs with 100% accuracy.

### 3.2. Host classification results

Identifying the origin of viral strains as human avian or swine has been done by scoring the pre-identified HA subtype with the corresponding “HA-Host” HMM models for better matching. “HMMSearch” program in HMMER suite has been used for this classification. The test results details of host classification for different HA subtypes in terms of accuracy, sensitivity and specificity are summarized in Table 2.

### 3.3. Model evaluation using ROC analysis

In a receiver operating characteristic (ROC) curve the true positive rate (sensitivity) is plotted in function of the false positive rate (100-specificity) for different cut-off points. Each point on the ROC plot represents a sensitivity/specificity pair corresponding to a particular decision threshold.<sup>27</sup> The following ROC curves were drawn using MedCalc program.<sup>28</sup> The curves indicate the observed criterion (threshold) values that maximized both sensitivity and specificity values. ROC curves for host identification of different HA subtypes are indicated in Figs 1–5.

## 4. Discussion

The obtained results confirm that profile HMM can successfully be used for classifying all influenza A stains hosted in all species in two major steps. First through

Influenza A Subtyping and Host Origin Classification

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43

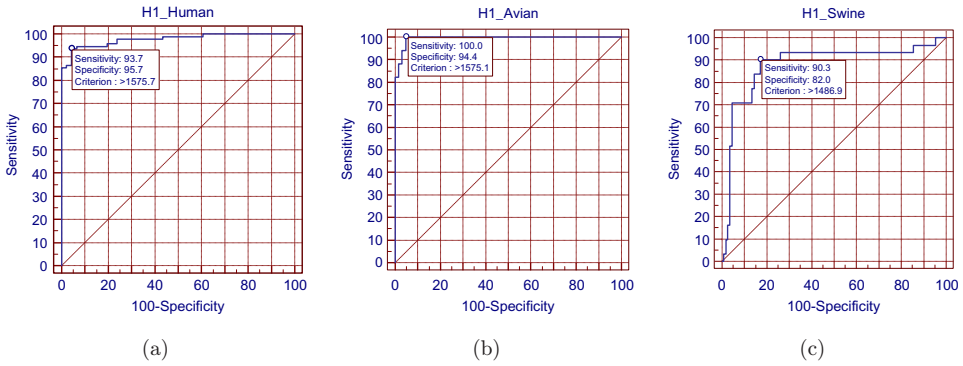


Fig. 1. ROC curves for host classification results of H1-human, H1-Avian and H1-Swine using HMM.

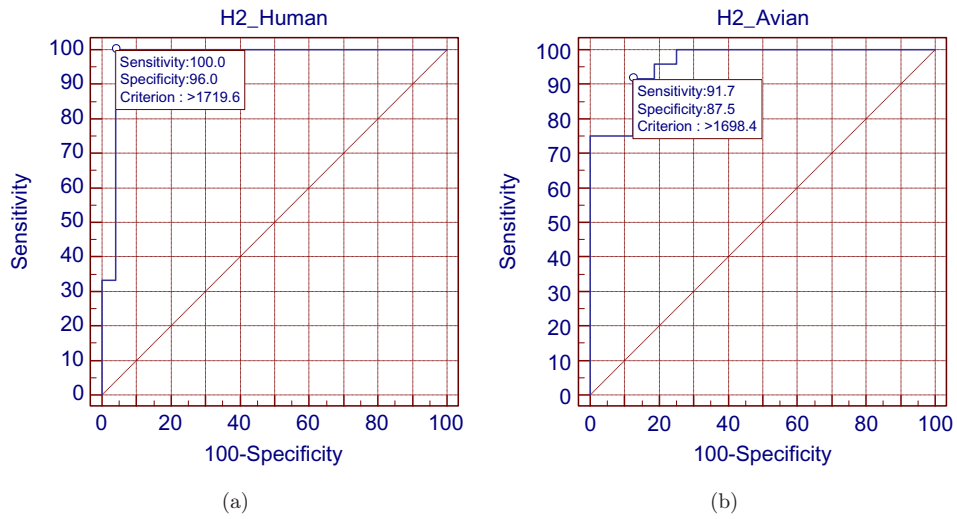


Fig. 2. ROC curves for host classification results of H2-human and H2-Avian using HMM.

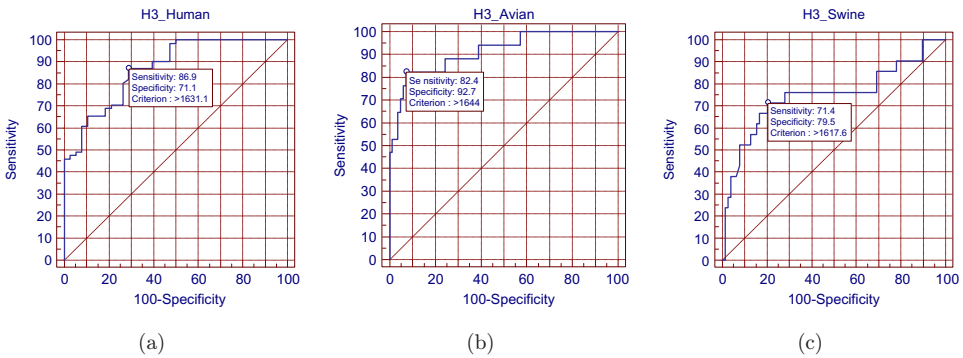


Fig. 3. ROC curves for host classification results of H3-human, H3-Avian and H3-Swine using HMM.

F. F. Sherif, M. El-Hefnawi & Y. M. Kadah

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43

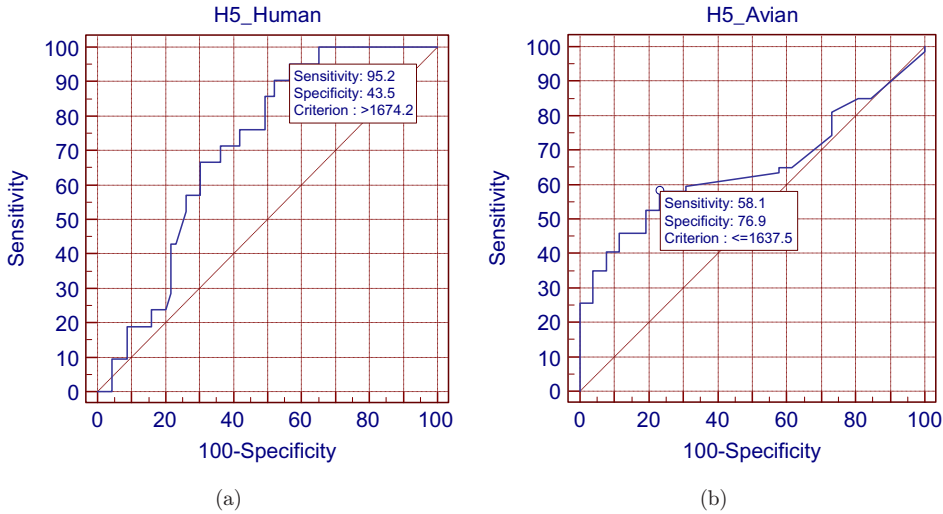


Fig. 4. ROC curves for host classification results of H5-human and H5-Avian using HMM.

identifying HA and NA subtypes. Second through predicting the host of origin of the pre-identified HA subtypes, by scoring it with the corresponding “HA subtype-Host HMM” models, searching for the best match.

For example, if a query HA sequence has been searched with each HA model separately, and we get the highest score with H1 model for example, then the entire sequence will be further scored with each H1 specified host separately: H1-human, H1-Swine and H1-Avian models searching for the highest match.

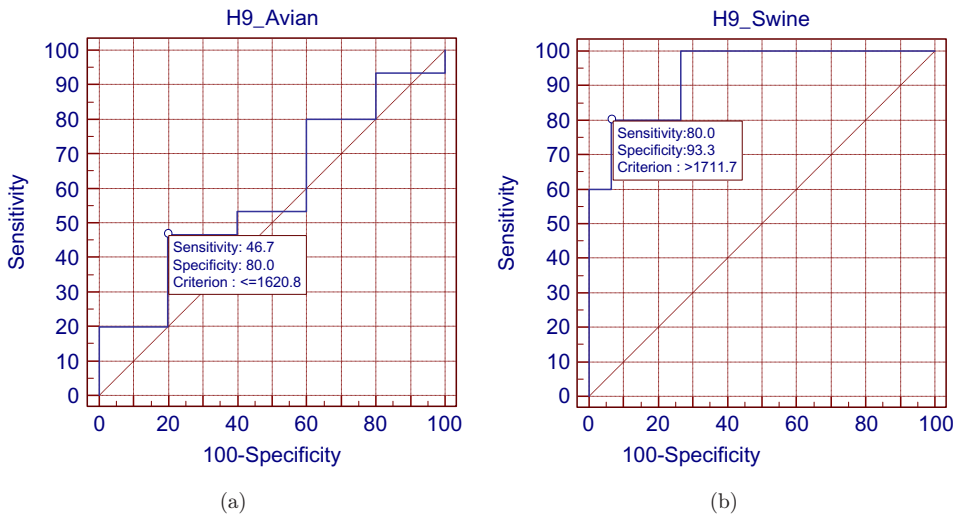


Fig. 5. ROC curves for host classification results of H9-Avian and H9-Swine using HMM.



1 All the 16 HA and 9 NA models have sensitivity of 100% and specificity of 100%.  
2 Although, there are differences in the criteria used in Attaluri *et al.*'s study and this  
3 study, their findings may support our findings that any unknown viral strain of  
4 influenza A, can be easily distinguishable as they have an extensive genetic diversity  
5 in HA and NA subtypes. Notably, our results achieved higher accuracy over Attaluri  
6 *et al.*'s study.<sup>13</sup> On the other hand, host classification of any viral sequence as human,  
7 avian or swine varied according to HA subtype. Among HA subtypes, there were  
8 some HAs (H1, H2, H3, H5 and H9) that can infect more than one species, through  
9 transmission of the whole virus or ever, the reassortment between avian and human  
10 viruses. Also, we found that some of those HA subtypes which can infect more than  
11 one species; vary greatly between human, swine and avian viruses. While some others  
12 vary little so it was difficult to identify their host of origin.

13 By comparing our results, we found that, H2 HA models have a higher accuracy  
14 over H1, H3, H5 and H9 HAs models. These results indicate that H2 viral subtypes  
15 have more genetic diversity between human and avian, compared to the other  
16 subtypes. In contrast, H5 HA models accuracies were not much higher than 53% for  
17 H5-Human and 63% for H5-Avian. This means that, no significant differences can be  
18 detected between human and avian H5 viruses using HMM.

19 These results agree with previous findings in Refs. 29 and 30, that highly  
20 pathogenic avian influenza H5N1 virus strains can transmit directly from avian  
21 species to humans and cause severe disease. The receptor binding preference of H5N1  
22 viruses can be altered by only a few amino acid substitutions in the HA protein. H1  
23 HA has accuracies of 94.4%, 84.5% and 89.5% for H1-human, H1- Swine and H1-  
24 Avian models, respectively. The host classification of H3 HA has the accuracies of  
25 80.8%, 78.7% and 90.9% for H3-human, H3-Swine and H3-Avian models, respec-  
26 tively. These results seem reasonable as cross-species infections usually take place in  
27 these subtypes, through reassortment or through whole host shift events. Never-  
28 theless, further improvement may be required in host classification to achieve higher  
29 accuracy. The remaining subtypes of HA are found only in avian hosts, so once they  
30 are classified by their subtypes as H4, H6, H7, H8, H10-H16, etc. they are also  
31 identified as having an avian host specification.

## 32 33 **5. Conclusions**

34  
35 Accurate detection of influenza viral origin and subtyping can significantly improve  
36 influenza surveillance and vaccine development. In this study, host identification and  
37 subtyping of influenza A virus were done based on HMMs for each subtype and major  
38 hosts (humans, avian, and swine). This study demonstrated the power of integrating  
39 the multiple sequence alignment and profile HMM approaches in classifying influenza  
40 A viral stains and their host of origin. In conclusion, our results indicate that influenza  
41 A sequences are HA and NA subtype specific and highly sensitive against HMM models  
42 (H1-H16), (N1-N9) and can easily be predicted with 100% accuracy. Host classification  
43 has accuracies that vary between 53% and 95.1% according to HA subtype.

F. F. Sherif, M. El-Hefnawi & Y. M. Kadah

## References

1. Horimoto T, Kawaoka Y, Influenza: Lessons from past pandemics, warnings from current incidents, *Nat Rev Microbiol* **3**:591–600, 2005.
2. Triki H, Clinical virology laboratory, *Arch Inst Pasteur Tunis* **74**:51–55, 1997.
3. Wiley DC, Skehel JJ, The structure and function of the hemagglutinin membrane glycoprotein of influenza virus, *Annu Rev Biochem* **56**:365–394, 1987.
4. Chander Y, Jindal N, Stallknecht DE, Sreevatsan S, Goyal SM, Full length sequencing of all nine subtypes of the neuraminidase gene of influenza A viruses using subtype specific primer sets, *J Virol Methods* **165**:116–120, 2010.
5. Munch M, Nielsen LP, Handberg KJ, Jorgensen PH, Detection and subtyping (H5 and H7) of avian type A influenza virus by reverse transcription-PCR and PCR-ELISA, *Arch Virol* **146**:87–97, 2001.
6. Zhang Y, Lin X, Zhang F, Wu J, Tan W, Bi S, Zhou J, Shu Y, Wang Y, Hemagglutinin and neuraminidase matching patterns of two influenza A virus strains related to the 1918 and 2009 global pandemics, *Biochem Biophys Res Commun* **387**:405–408, 2009.
7. Garten RJ, Davis CT, Russell CA, Shu B, Lindstrom S, Balish A, Sessions WM, Xu X, Skepner E, Deyde V, Okomo-Adhiambo M, Finelli L, Bridges CB, Shaw M, Jernigan DB, Uyeki TM, Smith DJ, Klimov AI, Cox NJ, Antigenic and genetic characteristics of swine-origin 2009 A(H1N1) influenza viruses circulating in humans, *Science* **325**:197–201, 2009.
8. Pedersen JC, Hemagglutination-inhibition test for avian influenza virus subtype identification and the detection and quantitation of serum antibodies to the avian influenza virus, *Methods Mol Biol* **436**:53–66, 2008.
9. Starick E, Romer-Oberdorfer A, Werner O, Type- and subtype-specific RT-PCR assays for avian influenza A viruses (AIV), *J Vet Med B Infect Dis Vet Public Health* **47**:295–301, 2000.
10. Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ, Gapped BLAST and PSI-BLAST: A new generation of protein database search programs, *Nucleic Acids Res* **25**:3389–3402, 1997.
11. Lu G, Jiang L, Helikar RM, Rowley TW, Zhang L, Chen X, Moriyama EN, GenomeBlast: A web tool for small genome comparison, *BMC Bioinformatics* **7**(4):18, 2006.
12. Eddy SR, Profile hidden Markov models, *Bioinformatics* **14**:755–763, 1998.
13. Attaluri PK, Chen Z, Weerakoon AM, Lu G, Integrating decision tree and hidden Markov model (HMM) for subtype prediction of human influenza A virus, in *Cutting-Edge Research Topics on Multiple Criteria Decision Making*, Springer, 2009, pp. 52–58.
14. Attaluri PK, Zheng X, Chen Z, Lu G, Applying machine learning techniques to classify H1N1 viral strains occurring in 2009 flu pandemic, in *The 6th Annual Biotechnology and Bioinformatics Symposium (BIOT-2009)*, The University of Nebraska-Lincoln, Lincoln, Nebraska, 2009.
15. Attaluri PK, Chen Z, Lu G, Applying neural networks to classify influenza virus antigenic types and hosts, in *2010 IEEE Symp Computational Intelligence in Bioinformatics and Computational Biology (CIBCB)*, Montreal, Canada, May 2–5, 2010.
16. Bao Y, Bolotov P, Dernovoy D, Kiryutin B, Zaslavsky L, Tatusova T, Ostell J, Lipman D, The influenza virus resource at the National Center for Biotechnology Information, *J Virol* **82**:596–601, 2008.
17. Suzuki Y and Nei M, Origin and evolution of influenza virus hemagglutinin genes, *Mol Biol Evol* **19**:501–509, 2002.
18. Chenna R, Sugawara H, Koike T, Lopez R, Gibson TJ, Higgins DG, Thompson JD, Multiple sequence alignment with the Clustal series of programs, *Nucleic Acids Res* **31**:3497–3500, 2003.

*Influenza A Subtyping and Host Origin Classification*

- 1 19. Jeanmougin F, Thompson JD, Gouy M, Higgins DG, Gibson TJ, Multiple sequence  
2 alignment with Clustal X, *Trends Biochem Sci* **23**:403–405, 1998.
- 3 20. Larkin MA, Blackshields G, Brown NP, Chenna R, McGettigan PA, McWilliam H,  
4 Valentin F, Wallace IM, Wilm A, Lopez R, Thompson JD, Gibson TJ, Higgins DG,  
5 Clustal W and Clustal X version 2.0, *Bioinformatics* **23**:2947–2948, 2007.
- 6 21. Schuster-Bockler B and Bateman A, An introduction to hidden Markov models, *Curr*  
7 *Protoc Bioinformatics*, Appendix 3, p. Appendix 3A, 2007.
- 8 22. Hughey R, Krogh A, Hidden Markov models for sequence analysis: Extension and anal-  
9 ysis of the basic method, *Comput Appl Biosci* **12**:95–107, 1996.
- 10 23. Finn RD, Tate J, Mistry J, Coghill PC, Sammut SJ, Hotz HR, Ceric G, Forslund K, Eddy  
11 SR, Sonnhammer EL, Bateman A, The Pfam protein families database, *Nucleic Acids*  
12 *Res* **36**:D281–D288, 2008.
- 13 24. Wilson D, Madera M, Vogel C, Chothia C, Gough J, The SUPERFAMILY database in  
14 2007: Families and functions, *Nucleic Acids Res* **35**:D308–D313, 2007.
- 15 25. Krogh A, Brown M, Mian IS, Sjolander K, Haussler D, Hidden Markov models in com-  
16 putational biology. Applications to protein modeling, *J Mol Biol* **235**:1501–1531, 1994.
- 17 26. Barrett C, Hughey R, Karplus K, Scoring hidden Markov models, *Comput Appl Biosci*  
18 **13**:191–199, 1997.
- 19 27. Nomura Y, Significance of the ROC (receiver operating characteristics) curve in diag-  
20 nostic tests, *Nippon Rinsho Jpn J Clin Med* 1402–1404, 1979.
- 21 28. Mariakerke B, MedCalc Software for Windows, Version 11.3.8 ed, 1993–2010
- 22 29. Wong SS, Yuen KY, Avian influenza virus infections in humans, *Chest* **129**:156–168,  
23 2006.
- 24 30. Auewarakul P, Suptawiwat O, Kongchanagul A, Sangma C, Suzuki Y, Ungchusak K,  
25 Louisirootchanaikul S, Lerdsamran H, Pooruk P, Thitithanyanont A, Pittayawonganon  
26 C, Guo CT, Hiramatsu H, Jampangern W, Chunsutthiwat S, Puthavathana P, An avian  
27 influenza H5N1 virus that binds to a human-type receptor, *J Virol* **81**:9950–9955, 2007.  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43