# INFLUENZA DIFFERENTIATION AND EVOLUTION\*

Krzysztof Bartoszek

Mathematical Statistics Chalmers University of Technology and University of Gothenburg Gothenburg, Sweden krzbar@chalmers.se

Pietro Liò<sup>†</sup>

Computer Laboratory, University of Cambridge Cambridge, United Kingdom pl219@cam.ac.uk

Anil Sorathiya

Computer Laboratory, University of Cambridge Cambridge, United Kingdom as833@cam.ac.uk

(Received February 18, 2010)

The aim of the study is to do a very wide analysis of HA, NA and M influenza gene segments to find short nucleotide regions, which differentiate between strains (*i.e.* H1, H2, ... *etc.*), hosts, geographic regions, time when sequence was found and combination of time and region using a simple methodology. Finding regions differentiating between strains has as its goal the construction of a Luminex microarray which will allow quick and efficient strain recognition. Discovery for the other splitting factors could shed light on structures significant for host specificity and on the history of influenza evolution. A large number of places in the HA, NA and M gene segments were found that can differentiate between hosts, regions, time and combination of time and region. Also very good differentiation between different Hx strains can be seen. We link one of our findings to a proposed stochastic model of creation of viral phylogenetic trees.

PACS numbers: 87.23.Kg, 87.10.Vg, 87.18.Tt

<sup>\*</sup> Presented at the Summer Solstice 2009 International Conference on Discrete Models of Complex Systems, Gdańsk, Poland, June 22–24, 2009.

<sup>&</sup>lt;sup>†</sup> Corresponding author.

### 1. Introduction

Statistical studies have provided interesting insights into the past influenza epidemics and pandemic. Every year the spreading of seasonal influenza causes significant mortality, and it cost billions of euros in health care expenses and with the risk of hospital-acquired infections [1–3]. Should the strains of influenza mutate to a strain that can quickly pass between birds and humans and then humans-humans it could cause a pandemic as it could spread to different parts of the world through human travel, migratory birds and in each place rapidly between humans.

The influenza virus is a common cause of respiratory infection all over the world. It infects not only humans, but also other species including avian and swine. Influenza A subtype (H1N1) could cause the next pandemic, if a new influenza A subtype has the ability to spread between humans efficiently. We have performed bioinformatic analysis to investigate the evolution of the HA, M and NA gene among different species. It would be important to know specific combination of viral RNA segments, hosts, regions, time and combination of time and region which have significant for host specificity and on the history of influenza evolution. Here we present a statistical analysis of influenza, and discuss how these statistics can provide insights on structures significant for host specificity and on the history of influenza evolution.

### 2. Data

Sequences of influenza strains were downloaded from GISAID [4]. Sequences from the HA, M and NA gene segments were downloaded. Each gene segment was considered separately. As some HA sequences were extremely short compared to the others these were removed. Each gene segment was downloaded with related information to it (*e.g.* its time period, geographical region and host). Unfortunately, for the sequences from the HA gene segment we did not have any information other than their strains. A summary of the downloaded gene segment sequences is presented in Table I.

In Table I the numbers of sequences do not always add up, there are two reasons for this. One is that for some sequences there is missing data *e.g.* host information but this was negligible. The main reason was that there were also sequences from the Australia and Oceania region and these were not considered in the regional analysis. The motivation was that there were relatively few of them compared to the other regions and when they were included nothing interesting was seen for them. Also a word is needed to motivate the method of regional splitting and temporal splitting. The temporal splitting was done in a fashion to have a equal balance between time bins. Unfortunately, the further we go back in time the fewer and fewer sequences we have. This type of binning has the consequence of ignoring

#### TABLE I

Breakdown of downloaded sequences. Abbreviations are: EA — Eastern Asia, CA — Central Asia, NA — North America, AO — Australia and Oceania, MEI — Middle East,Indian Peninsula, WE — Western Europe, SA — South America, EES — Eastern Europe Scandinavia, UKII — UK, Ireland, Iceland, A — Africa, EuA — Eurasia, As — Americas, Av — Avian, Sw — Swine, H — Human, Eq — Equine, Env — Environment, Om — Other mammals, Un — Unknown (these abbreviations will be used in subsequent tables).

HA gene segm	ent								
	6052 sequences		alignme	nt length: 20	63 ba	ases			
H1	H2	H3	H4	H5	H6	H7	H8		
904	132	1608	158	1833	286	509	17		
H9	H10	H11	H12	H13	H14	H15	H16		
387	65	71	26	27	4	7	15		
HA gene segm	ent H1N1								
	1851  seque	ences		alignment	leng	th: 2	2063 1	oases	
EA	CA	NA	AO	MEI	WE	$\mathbf{SA}$	EES	UKII	Α
217	51	980	192	23	64	83	96	107	22
Human	Swine	Avian	Lab						
1570	191	77	10						
M gene segme	nt								
6775 sec	quences	alignmer	nt length:	1101 bases					
EuA+Africa	As	1902 - 1999	2000-200	4 2005-2009					
2670	3243	2051	2227	2494					
EuA+Africa	EuA+Africa	EuA+Africa							
1902 - 1999	2000 - 2004	2005 - 2009							
616	956	1053							
As $1902 - 1999$	As $2000-2004$	As $2005-2009$							
1346	650	1247							
Av	Sw	Н	Eq	Env	Om	Un			
2578	320	3650	90	101	18	18			
NA gene segm	ent								
9661 see	quences	alignmer	nt length:	1741  bases					
EuA	As	1902 - 1999	2000-200	4 2005–2009					
4475	4160	2433	3244	3984					
EuA 1902–1999	EuA 2000–2004	EuA 2005–2009	)						
815	1729	1931							
As 1902–1999	As 2000–2004	As 2005–2009							
1502	842	1816							
Av	Sw	Н	Eq	Env	Om	Un			
3827	461	5121	89	102	30	19			

any changes that happened inside the time bins. The main motivation for geographic splitting was land and water masses. The whole of Eurasia (and effectively with Africa) is connected by land while it is separated from the Americas (the South and North also connected by land) and Australia and Oceania by water. These divisions of time and space are very rough and are such to keep the number of bins at a low level.

# 3. Method

The method employed to analyze a set of sequences is a very preliminary one. The aim is to have a first look at the data and see whether anything new can be discovered without explicit mathematical modelling. For computational reasons mathematical models (*e.g.* Markov models) of alignment data nearly always assume independence of columns. We are not aware of any satisfactory and computationally effective model that considers even short dependencies between columns in such a large alignment. Therefore we want to see how much can be discovered without the additional knowledge a model would bring. In this paper we will introduce the following terminology, a metasubsequence of length k will be vector of k consecutive positions in the alignment and by the value of the metasubsequence we will mean the actual values (for a given sequence, A, C, G, T in a nucleotide alignment) that are observed in the respective columns of the alignment.

# 3.1. Alignment

The alignment of the sequences was done in Clustal 2.0.10 [5]. Due to the large number of sequence it was done with the option approximate and all others default. Prior to the alignment sequences which were considerably shorter than the others had to be removed. In the case of the H1N1 sequences this was roughly 1500 sequences. Had this not been done the contingency table test did not look at splits between hosts, regions, *etc.* but between short and long sequences.

### 3.2. Contingency table

The contingency table test (or chi-square test) is described in detail in e.g. [6] but we will make a short overview of it here. Let us assume we have N objects in a two-way categorization table with an arbitrary number of rows and columns. In our case the rows will be the different possible values of metasubsequences the regions take and the columns will describe the different possible categories the sequences can take). Such a table can look like the one in [6].

					Total
	$Y_{11}$	$Y_{12}$		$Y_{1k}$	$y_{1\cdot}$
	$Y_{21}$	$Y_{22}$		$Y_{2k}$	$y_2$ .
	:	÷	·	÷	÷
	$Y_{n1}$	$Y_{n2}$		$Y_{nk}$	$y_n$ .
Total	$y_{\cdot 1}$	$y_{\cdot 2}$		$y_{\cdot k}$	y

Under the null hypothesis that there is no association between row and column categories each entry  $Y_{jk}$  will be a random variable with mean value of  $E_{jk} = \frac{y_j \cdot y_k}{y}$ . The test statistic is  $\sum_{jk} \frac{(Y_{jk} - E_{jk})^2}{E_{jk}}$  and under the null hypothesis has a chi-square with (r-1)(c-1) degrees of freedom distribution. An important assumption of this test is that the observations leading up to the counts are independent. In our case they are not due to the phylogenetic history. Taking this into account is a topic for further study and due to this we cannot assign any formal statistical significance to the results. This problem of dependence has been very recently looked in a simulation study in [7] however the authors consider alignment columns independently. The procedure of working with the contingency table is described in the following algorithm

- 1: for metasubsequence length i = 0 to n do
- 2: for column j = 1 to length of alignment do
- 3: build contingency table for metasubsequence at position j to j + i{i is from 0}
- 4: calculate *p*-value of contingency table
- 5: end for
- 6: end for
- 7: return those contingency tables with their positions and metasubsequences that are below some cut-off p-value
- 8: go through the returned contingency tables and see whether they are interesting.

The cut-off *p*-value has to be extremely low (it was taken to be  $1^{-183}$ ) due to the fact that we have at each position very often one sequence which is totally different from the others and such a sequence could immediately generate a significant table. Unfortunately very often we have a huge number of tables, in fact in some cases every single consecutive metasubsequence of the alignment was significant due to this effect. Our choice of scoring is presented in the next section. Up until now most analysis have been treating positions independently *e.g.* [8] (nucleotide alignment) [9] (protein alignment). An independent site analysis by a direct application of entropy [9–11] methods is straightforward, quick and the results are immediate and do not require any postprocessing. The problem is that it might miss columns that are significant only when a combination of them is considered. Here we make an attempt to do a large scale sequence analysis to find metasubsequences of up to 25 bases which categorize the sequences. We are especially interested in metasubsequences of up to 25 bases because this is what we expect that the length of the probe on the Luminex microarray will be.

#### 3.3. Scoring contingency tables

The difficulty with the contingency table is that if the alignment is very "noisy" (either due to misalignment or a large number of sequences that have mutations) nearly all positions in the alignment will turn significant. Therefore some other method than the *p*-value of scoring each contingency table has to be used. We adopted the following strategy for each contingency table,

1: remove all rows that have in total less than p% of the sequences

- 2: remove all columns that have in total less than p% of the sequences
- 3: N = number of rows
- 4: for each column i do

$$Score_i = 1 - \frac{-\sum_{\text{each row j}} p_{ij} \log p_{ij}}{\log N}$$

#### 6: end for

7: return  $\max_i$  Score<sub>i</sub>.

We look at each column of the contingency table (which represents some grouping of the sequences) and see what is the estimated entropy of the sequences in this group  $-\sum_{\text{each row } i} p_{ij} \log p_{ij}$  dividing by  $\log N$  the maximal entropy possible. The score will be 1 if we can predict the value of the metasubsequence perfectly, *i.e.* there is exactly one value of a given metasubsequence in the given group and 0 if we cannot say anything. To get the score for a given position (in Figs. 2, 3, 4 and 5) we take the maximum over all metasubsequences of length  $1 \dots 25$ . In our analysis we decided to ignore gaps in the HA and M gene segment alignments as these were aligned with very few gaps (except on the edges) while in the analysis of the NA gene segment gaps were treated as a fifth residue (here there was a large number of gaps inside the alignment). The methodology presented is naturally very heuristic due to the large computational intensity and complexity of the problem. However, it guarantees that if a given metasubsequence is present in enough sequences and its values split perfectly between groups then the method will find it.

### 3.4. Implementation

The described method was implemented as a Perl [12] script. For working with the alignments it uses BioPerl [13] and for calculating the *p*-value of the  $\chi^2$  statistic the module Statistics::Distributions [14] was used.

### 3.5. Phylogenetic tree

HA gene segment of 6052 sequences was aligned using Clustal 2.0.10 tool [5]. This alignment of 2063 bases length from H1 to H16 viral strains was used to generate the phylogenetic tree which describing the evolution of strains from one to another (see figure 1). Colour coding of the parenthesis in Fig. 1 represents strains from H1 to H16 and Dendroscope 2.4 tool [15] was used to visualize the phylogenetic tree.



Fig. 1. Phylogeny tree of influenza strains; Each colour represents each strain of influenza virus.

### 4. Results

The biological questions posed in the analysis are the following (each item is done on a different sequence set),

• Can we find regions of length 20–25 bases the distinguish between Hx strains (for the Luminex microarray)?

- Can we find short regions up to 25 bases which will distinguish between hosts and regions in H1N1 sequences?
- Can we find short regions up to 25 bases which will distinguish between hosts, regions, time periods in M sequences, can any temporal and geographical interactions be found?
- Can we find short regions up to 25 bases which will distinguish between hosts, regions, time periods in NA sequences, can any temporal and geographical interactions be found?

In all cases interesting places were found.

# 4.1. Hx analysis

Only H1, H2, H3, H4, H5, H6 and H7 sequences were considered. The other ones were too few in numbers compared to these, so were discarded in the analysis as they only acted as "noise". All in all 292 places of length 20–25 bases were found in the HA genomes that split (nearly perfectly) between the strains. Some of them built up together to form a larger region so we could actually distinguish about 6 regions in the genome where differences between the strains can be found. It is not possible to present all the results but a nice example contingency table can be shown, Table II. The values inside the table are the fraction of the number of strains having the given value of the metasubsequence. The fact that for a strain a column does not add up to 1 is due to it having less than 100 sequences in some rows and for clarity of the presentation if a row contained less than 100 sequences it was not written out. The scoring function failed to cut-down on the number of positions as all of the splits were perfect. In Fig. 2 we can see how that the scores were nearly 1 all along the gene segment.

### TABLE II

Example contingency table splitting Hx strains of HA gene segment, 292 such places were found. We can see that combinations of different positions are needed to differentiate between all.

Position 1687	H1	H2	H3	H4	H5	H6	H7
TGGGACTTATGACCATGATGTATAC	0	0	0.07	0	0	0	0
CAACACTTATGACCATACTCAATAC	0	0	0	0	0	0	0.24
CAACACGTATGACCATACTCAATAC	0	0	0	0	0	0	0.21
TGGAACTTATGACCATGATGTATAC	0	0	0.53	0	0	0	0
TGGAACTTATGACCACGATGTATAC	0	0	0.1	0	0	0	0
TGGAACTTATGACTATCCAAAATAT	0.63	0	0	0	0	0	0
$\mathbf{CGGA}\mathbf{A}\mathbf{C}\mathbf{G}\mathbf{T}\mathbf{A}\mathbf{T}\mathbf{G}\mathbf{A}\mathbf{C}\mathbf{T}\mathbf{A}\mathbf{C}\mathbf{C}\mathbf{C}\mathbf{G}\mathbf{C}\mathbf{A}\mathbf{G}\mathbf{T}\mathbf{A}\mathbf{T}$	0	0	0	0	0.69	0	0



Fig. 2. Scores along the alignment of the HA gene segment. We can see that nearly all along the gene segment we have a score of 1. The visible gaps in the graph are due to gaps in the alignment in those places.

### 4.2. H1N1 analysis

Two types of analysis were done, whether any place specific for hosts or region can be found.

In the host analysis (Table III) five places specific for human hosts were found. They were at positions 282–307, 1085–1110, 1133–1158, 1219–1244 and 1570–1575 in the alignment. About 83% of all human had the same sequence.



Fig. 3. Top: scores along the alignment of the HA H1N1 gene segment for splits between hosts. All along the gene segment we have positions achieving high scores (cut-off for graph is 0.99). A few of the interesting examples are in Table III. Bottom: scores along the alignment of the HA H1N1 gene segment for splits between regions. All along the gene segment we have positions achieving high scores (cut-off for graph is 0.99). A few of the interesting examples are in Table IV.

As previously the columns do not add up to 1 due to many rows having very small numbers of sequences and we can see changes specific for all three hosts. When the analysis was done according to regions the North American region stood out significantly (Table IV) in the HA gene segment.

e Avian	1 0	e Avian	$\begin{array}{ccc} 8 & 0.052 \\ 8 & 0.3 \\ 0.13 \\ 0.13 \end{array}$
Swine	0.01	Swine	$\begin{array}{c} 0.408 \\ 0.183 \\ 0.026 \end{array}$
Human	0.836 0.161	Human	$\begin{array}{c} 0.002 \\ 0.001 \\ 0.98 \end{array}$
Position 513(1)	Ηď	Position 1128 (0.96)	GGTTTTATTGAGGG GGATTCATTGAAGG GGTTTCATTGAAGG
Avian	0.065 0.935	Avian	$0.961 \\ 0$
Swine	0.728 0.272	Swine	$0.984 \\ 0$
Human	0.004 0.996	Human	0.305 0.692
Position 396(0.96)	ЧÜ	Position 1052(0.98)	ΔA
Avian	$\begin{array}{c} 0.013\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0 \end{array}$	Avian	1 0
Swine	$\begin{array}{c} 0.387 \\ 0.021 \\ 0.01 \\ 0 \\ 0 \end{array}$	Swine	10
Human	$\begin{array}{c} 0.003\\ 0.042\\ 0.1\\ 0.734\\ 0.045 \end{array}$	Human	$0.175 \\ 0.824$
Position 286 (0.96)	TGGGAAACCCAGAATGTGAATTACT TGGGAAATCCAGAGTGTGAATCACT TAGGAAACCCAGAATGCGAATCACT TAGGAAACCCAGAATGCGAATCACT TAGGAAACCCAGAATGCGAATTACT TAGGAAACCCAGAATGCGAATTGCT	Position 579 (1)	ΥÜ

Avian	0	0	0	0	0	0.013	0	0	Avian	0	0.013	0	0	0
Swine	0	0	0	0.513	0.01	0	0	0	Swine	0.021	0	0	0	0.225
Human	0.03	0.03	0.024	0.042	0.036	0.682	0.045	0.055	Human	0.042	0.825	0.024	0.027	0
Position 1261 (0.97)	TGAACTCTATAATCGAGAAAATGAA	TGAATTCTGTGATTGAGAAAATGAA	TGAACTCTGTTATCGAGAAAATGAA	TAAATTCTGTTATTGAAAGATGAA	TGAATTCTGTAATCGAGAAAATGAA	TGAATTCTGTAATTGAGAAAATGAA	TGAATTCTGTAATTGAGAAGATGAA	TGAACTCTGTAATTGAGAAGATGAA	Position 1582 (0.96)	ACCCAAAATACTCAGAGGAAGCAAA	ATCCAAAATATTCCGAAGAATCAAA	ATCCAAAATATTCAGAGGAATCAAA	ATCCAAAATATTCAGAAGAATCAAA	ACCCAAAGTACTCAGAAGAATCAAA
Avian	0	0	0.974						Avian	0.182	0.818			
Swine	0	0.01	0.942						Swine	1	0			
Human	0.862	0.031	0.096						Human	1	0			
Position 1155 (0.98)	ATGGTAGATGG	ATGATTGATGG	ATGATAGATGG						Position 1566 (0.99)	AATGG	AACGG			

TABLE III

Examples of top scoring (above 0.95) host specific changes HA genome of H1N1, the score is in brackets next to the position

number. There were in total 70 such places found.

426

TABLE IV

4 0 0 0 -K 0 0 ----K 0 0 -0 4  $\triangleleft$ 00 0 -- $\begin{array}{c} 0.009 \\ 0.075 \\ 0.009 \end{array}$  $0.009 \\ 0.075$ 0.8790.8880.0190.8690.0280.8690.019UKII UKII UKII 0.028UKII UKII 0.8690.0840 0.9790.0520.927 $0.01 \\ 0.969$ 0.979EES 0.01EES 0.99EES 0.01EES EES 0.01 $\begin{array}{c} 0 \\ 0.01 \end{array}$ 0.01 $0 \\ 0.01$ 0 0.9880.8430.157SASA0 0  $\mathbf{S}\mathbf{A}$ SA $\mathbf{S}\mathbf{A}$ 0 1 0 0 0 1 0 1 0 $0.078 \\ 0.4212$ 0.0160.4060.4690.453 $0.234 \\ 0.469$ 0.0160.328WE 0.0310.531WE WE  $0.5 \\ 0.031$ WE WE 0.50 0.913 $0.826 \\ 0.174$ MEI MEI MEI MEI MEI 0 0 0 1 0 0 0 1 0 1 0 0.0160.0160.0680.9690.9840.0630.0260.9740.9220.911AO 0.01 AO AO AO 0.01 AO 0 0 0 0 0.2030.7320.0460.0140.9090.1880.7360.0450.733 $0.573 \\ 0.363$ 0.0590.0460.010.2510.01NA NA NA NA NA 0.0980.0780.0590.039 $0 \\ 0.098$ 0.137 0.7650.8820.7840.8630.1370.784CA CA CA CA CA 0 0 0 Score is in brackets next to position number 0.0920.8430.8760.0050.0140.0730.8990.0090.9720.1060.9820.0050.037EA EA EA EA EA 0 0 0 Position 649 (0.96) Position 333 (0.99) Position 286(1)Position 358(1)Position 507(1)TTGGCAA TGGGAAA ACAATGG AGAATGG TGTGGGG TCTGGGG TGGGCAA TAGGAAA TATATAA TACATTG AAATGG AGAACGG TATGGGG TACATAA GGAGT GGTGT

Examples of top scoring (above 0.95) regional differentiation of HA H1N1 sequences. There were 229 such positions found.

Position $706(1)$	EA	CA	NA	AO	MEI	WE	$\mathbf{SA}$	EES	UKII	A
AA CT CA	$\begin{array}{c} 0.866 \\ 0.009 \\ 0.101 \end{array}$	$\begin{array}{c} 0.804 \\ 0.039 \\ 0.078 \end{array}$	$\begin{array}{c} 0.742 \\ 0.044 \\ 0.194 \end{array}$	$\begin{array}{c} 0.974 \\ 0 \\ 0.016 \end{array}$	$\begin{array}{c} 0.957\\ 0\\ 0\\ 0 \end{array}$	$\begin{array}{c} 0.469 \\ 0.359 \\ 0.031 \end{array}$	1 0 0	$\begin{array}{c} 0.979\\0\\0\\0.01\end{array}$	$\begin{array}{c} 0.869 \\ 0.084 \\ 0.047 \end{array}$	0 0
Position $902(1)$	EA	CA	NA	AO	MEI	WE	$\mathbf{SA}$	EES	UKII	A
GGC AAT GGG GGT	$\begin{array}{c} 0.88 \\ 0 \\ 0.0046 \\ 0.074 \end{array}$	$\begin{array}{c} 0.88 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0.059 \end{array}$	$\begin{array}{c} 0.75 \\ 0.064 \\ 0.0561 \\ 0.120 \end{array}$	$\begin{array}{c} 0.98 \\ 0.016 \\ 0 \\ 0 \end{array}$	$\begin{array}{c}1.00\\0\\0\\0\end{array}$	$\begin{array}{c} 0.92 \\ 0.031 \\ 0.0312 \\ 0 \end{array}$	1000	$\begin{array}{c} 0.99\\ 0.01\\ 0\\ 0\\ 0\\ 0\end{array}$	$\begin{array}{c} 0.94 \\ 0 \\ 0.028 \\ 0.028 \end{array}$	$\begin{array}{c} 0.91 \\ 0 \\ 0 \\ 0 \\ 0 \end{array}$
Position $947(1)$	EA	CA	NA	AO	MEI	WE	$\mathbf{SA}$	EES	UKII	A
TGCAC TGCGA	0.0046 0	$\begin{array}{c} 0.0588 \\ 0 \end{array}$	0 0.04	00	00	$\begin{array}{c} 0.1250\\ 0\end{array}$	0 0	$\begin{array}{c} 0.0104 \\ 0 \end{array}$	00	0 0
TGCAA TGTAC	0.018	00	0.069	0.016	00	0.031 0.125	00	$0.01 \\ 0$	0.028 0.084	00
TGTGA TGTAA	$0.88 \\ 0.092$	$\overset{0.8}{0.137}$	0.77 $0.124$	0.98 $0$	$\frac{1}{0}$	0.72 $0$	$\frac{1}{0}$	0.98 0	0.87 0.019	$\frac{1}{0}$
Position $1422 (1)$	EA	CA	NA	AO	MEI	WE	SA	EES	UKII	A
AGGACT AGAACT	$0.866 \\ 0.115$	$0.824 \\ 0.137$	$0.761 \\ 0.229$	$0.917 \\ 0.078$	$0.957 \\ 0.044$	$0.797 \\ 0.047$	$\begin{array}{c} 1\\ 0 \end{array}$	$0.958 \\ 0.042$	$0.916 \\ 0.019$	$\begin{array}{c} 1\\ 0 \end{array}$
Position 1566 (0.95)	EA	CA	NA	AO	MEI	WE	$\mathbf{SA}$	EES	UKII	A
AATGGAACTTA AATGGCACATA AATGGGACTTA AACGGCACATA	$\begin{array}{c} 0.899 \\ 0.009 \\ 0.0783 \\ 0.005 \end{array}$	$\begin{array}{c} 0.765 \\ 0.098 \\ 0.137 \\ 0\end{array}$	$\begin{array}{c} 0.688 \\ 0.001 \\ 0.248 \\ 0.054 \end{array}$	$\begin{array}{c} 0.974 \\ 0 \\ 0.0156 \\ 0 \end{array}$	$\begin{array}{c} 1\\ 0\\ 0\end{array}$	$\begin{array}{c} 0.453 \\ 0.422 \\ 0.0313 \\ 0.078 \end{array}$	$\begin{array}{c} 0.988 \\ 0 \\ 0.0121 \\ 0 \end{array}$	$\begin{array}{c} 0.979 \\ 0 \\ 0.0104 \\ 0 \end{array}$	$\begin{array}{c} 0.907\\ 0.075\\ 0.009\\ 0\end{array}$	0.909 0 0 0
Position $1611(1)$	EA	CA	NA	AO	MEI	WE	SA	EES	UKII	A
AG AA	$0.419 \\ 0.581$	$0.451 \\ 0.529$	$0.011 \\ 0.986$	$0.182 \\ 0.818$	$0.044 \\ 0.957$	$\begin{array}{c} 0.031 \\ 0.969 \end{array}$	$0.048 \\ 0.952$	$0.073 \\ 0.927$	0	0

### 4.3. M results

The M gene segment sequences were analyzed to find differences between hosts, regions, time periods and geographical and temporal interactions. The results are presented in Fig. 4 and Tables V, VI, VII and VIII.



Fig. 4. Top left: scores along the alignment of the M gene segment for splits between hosts, cut-off for graph is 0.99. Top scoring examples are in Table VII. Top right: scores along the alignment of the M gene segment for splits between regions, cut-off for graph is 0.95. Top scoring examples are in Table V. Bottom left: scores along the alignment of the M gene segment for splits between time periods, cut-off for graph is 0.95. Top scoring examples are in Table VI. Bottom right: scores along the alignment of the M gene segment for splits between regions combined with time periods, cut-off for graph is 0.99. Top scoring examples are in Table VIII.

# TABLE V

Position 143 Score 0.99	Eurasia & Africa	Americas	Position 242 Score 0.98	Eurasia & Africa	Americas
CAGAGACT CAGAAACT	$0.709 \\ 0.278$	$\begin{array}{c} 0.994 \\ 0 \end{array}$	ATTTT ATGTT	$0.867 \\ 0.067$	$\begin{array}{c} 0.985\\ 0\end{array}$
Position 538 Score 0.99	Eurasia & Africa	Americas	Position 683 Score 0.95	Eurasia & Africa	Americas
TC GC	$\begin{array}{c} 0.9 \\ 0.1 \end{array}$	$0.999 \\ 0.001$	TTGC TCGC	$0.627 \\ 0.369$	$\begin{array}{c} 0.99 \\ 0.005 \end{array}$
Position 768 Score 0.96	Eurasia & Africa	Americas	Position 944 Score 0.99	Eurasia & Africa	Americas
CTTGAAAATTT ATTGAAAATTT	$0.826 \\ 0.149$	$\begin{array}{c} 0.977 \\ 0 \end{array}$	CCTTCTACGGCAGG CCTGCTACGGCAGG CCTTCTACGGAAGG	$0.059 \\ 0.071 \\ 0.814$	$\begin{array}{c} 0\\ 0\\ 0.984 \end{array}$

Top scoring (above 0.95) region specific changes of M gene segment.

# TABLE VI

Top scoring (above 0.95) time period specific changes of M gene segment.

Position 143 Score: 0.97	Until 1999	$\begin{array}{c} 2000-\\ 2004 \end{array}$	2005 - 2009	Position 242 Score: 0.96	Until 1999	$\begin{array}{c} 2000-\\ 2004 \end{array}$	2005 - 2009
CAGAGACT CAGAAACT	$0.987 \\ 0.001$	$0.920 \\ 0.072$	$0.761 \\ 0.232$	ATTTT ATGTT	$0.973 \\ 0.0005$	$0.972 \\ 0.002$	$\begin{array}{c} 0.861 \\ 0.07 \end{array}$
Position 326 Score: 0.97	Until 1999	$\begin{array}{c} 2000-\\ 2004 \end{array}$	2005 - 2009	Position 571 Score: 0.96	Until 1999	$\begin{array}{c} 2000-\\ 2004 \end{array}$	2005 - 2009
AATGG AACGG	$0.893 \\ 0.106$	$0.991 \\ 0.009$	$0.996 \\ 0.0036$	AC AT	$0.996 \\ 0.004$	$\begin{array}{c} 0.888\\ 0.106 \end{array}$	$0.765 \\ 0.235$
Position 871 Score: 0.95	Until 1999	$\begin{array}{c} 2000-\\ 2004 \end{array}$	2005 - 2009	Position 939 Score: 0.95	Until 1999	$\begin{array}{c} 2000-\\ 2004 \end{array}$	2005 - 2009
GATATTG GATACTG	0.982 0.003	$0.899 \\ 0.093$	$0.966 \\ 0.0008$	GAGGCCCTTCTAC GAGGGCCTTCTAC GAGGGCCTGCTAC	$0.001 \\ 0.956 \\ 0$	$0.146 \\ 0.764 \\ 0.042$	$\begin{array}{c} 0.233 \\ 0.676 \\ 0.031 \end{array}$
Position 1023 Score: 0.98	Until 1999	2000– 2004	2005 - 2009				
GTCAT ATCAT	$0.992 \\ 0.001$	$0.994 \\ 0.002$	$0.827 \\ 0.164$				

TABLE VII

Equine

Human

Swine

Avian

Position 914(0.98)

Equine

Human

Swine

Avian

Position 478 (0.96)

Equine

Human

Swine

Avian

 $\begin{array}{c} \text{Position } 410 \\ (097) \end{array}$ 

Most interesting of the top scoring (above 0.95) host specific changes of M gene segment. Score in brackets after position number. There were 73 such positions found.

97 0.8 22 0 39 0 39 0	Equine	$\begin{array}{c} 0.944 \\ 0.056 \\ 0 \end{array}$	Equine	$\begin{array}{c} 0.889 \\ 0 \\ 0 \end{array}$
$\begin{array}{c} 0.85 & 0.0 \\ 0 & 0.12 \\ 0.003 & 0.0 \\ 0.003 & 0.0 \end{array}$	Human ]	$\begin{array}{c} 0.520\\ 0.207\\ 0.27\end{array}$	Human l	$\begin{array}{c} 0.095 \\ 0.281 \\ 0.589 \end{array}$
0.004	Swine	0.875 0.009 0.019	Swine	0.869 ).009375 0.059
CGTCG CGTAT CGATT CGATT	Avian	0.989 0.006 0	Avian	$\begin{array}{c} 0.952 \\ 0.0008 \\ 0 \end{array}$
-1 0	933	V V V	1013	GAT GAT GAC
0.082	osition (0.96)	TGAA TAAA TTAA	(0.95)	TGAC TGAC TGAC
0.928	- <sup>D</sup>		Pc	500 500
0.991	Equine	$\begin{array}{c} 0 \\ 0.011 \\ 0 \\ 0.9 \\ 0.9 \end{array}$	Equine	$\begin{array}{c} 0.911\\ 0\\ 0\\ 0 \end{array}$
AC GC	Human	$\begin{array}{c} 0.523 \\ 0.25 \\ 0.042 \\ 0.057 \\ 0.095 \end{array}$	Human	$\begin{array}{c} 0.724 \\ 0.01 \\ 0.243 \end{array}$
1 0	Swine	$\begin{array}{c} 0.0406\\ 0\\ 0.013\\ 0\\ 0.872\end{array}$	Swine	$\begin{array}{c} 0.947\\ 0\\ 0\end{array}$
0.907	Avian	$\begin{array}{c} 0\\ 0.00233\\ 0\\ 0\\ 0.976 \end{array}$	Avian	$\begin{array}{c} 0.888 \\ 0.054 \\ 0.009 \end{array}$
0.091		ជំពុំពិណ័ណ		- 
0.003	ion 921 0.98)	ACACG ACACG ACACG ACACG ACACG ATACG	ion 939 .98)	10001 10001
AT GT	Positi (0	TCAA, TTAA, TTGA, TCAG, TCAG,	Positi (0	GAGG GAGG GAGG

Influenza Differentiation and Evolution

Top scoring (above 0.9 African sequences not	5) regi includ	ons and ed in I	d time Eurasia	period	l specil	ic char	iges of M gene segment.	Score	in bra	ckets a	fter pos	ition n	umber.
Position 98 (0.97)	EuA u 99	${ m EuA}_{00-04}$	$_{05-09}^{\rm EuA}$	$^{\rm As}_{ m u \ 99}$	$^{\mathrm{As}}_{00-04}$	$^{\mathrm{As}}_{05-09}$	Position 143 (0.96)	EuA u 99	EuA 00-04	EuA $05-09$	As u 99	${}^{ m As}_{ m 00-04}$	$^{ m As}_{05-09}$
GTTCTCTCTAT GTTCTTTCTAT	0.130 0.039	0.937 0.053	$0.904 \\ 0.069$	$0.981 \\ 0.001$	0.983 0.002	0.886 0.112	CAGAAACTTGAAGATGT CAGAGACTTGAAGATGT CAGAGACTGGAGGATGT CAGAGACTGGAAAGTGT CAGAGACTGGAAAGTGT CAGAGACTTGAGGATGT	$\begin{array}{c} 0.002 \\ 0.818 \\ 0 \\ 0 \\ 0 \\ 0.065 \end{array}$	$\begin{array}{c} 0.166\\ 0.659\\ 0.002\\ 0\\ 0\\ 0.103 \end{array}$	$\begin{array}{c} 0.393\\ 0.310\\ 0.12\\ 0.038\\ 0.087\end{array}$	0 0.939 0 0.039	$\begin{array}{c} 0 \\ 0.971 \\ 0 \\ 0 \\ 0.009 \end{array}$	$\begin{array}{c} 0 \\ 0.794 \\ 0 \\ 0.0111 \\ 0.078 \end{array}$
Position 185 (0.95) GAGGCACT GAGGCTCT	EuA u 99 0.044 0.920	EuA 00 $-04$ 0.006 0.983	$EuA \\ 05-09 \\ 0.002 \\ 0.979$	As u 99 0.449 0.482	${\rm As} \\ 00-04 \\ 0.392 \\ 0.591$	$\substack{\text{As}\\05-09\\0.232\\0.711}$	Position 200 TGGCTAAAG TGGTTAAAG	EuA u 99 0.992 0.005	$EuA \\ 00-04 \\ 0.951 \\ 0.015$	$EuA \\ 05-09 \\ 0.93 \\ 0.006$	As u 99 0.006	$\substack{ {\rm As} \\ 00-04 \\ 0.992 \\ 0.008 \\ \end{array}$	$\substack{{\rm As}\\05-09\\0.92\\0.073}$
Position $204$ (0.97)	EuA u 99	${ m EuA}_{00-04}$	$_{05-09}^{\rm EuA}$	As u 99	$^{ m As}_{ m 00-04}$	$^{\mathrm{As}}_{05-09}$	Position 221 (0.95)	EuA u 99	EuA 00-04	EuA $05-09$	As u 99	${}^{ m As}_{ m 00-04}$	$^{ m As}_{05-09}$
TAAAGACAAGACCAAT TAAAGACAAGACCGAT	$0.929 \\ 0.06$	$0.976 \\ 0.017$	$0.956 \\ 0.01$	$0.906 \\ 0.085$	$0.942 \\ 0.051$	$0.993 \\ 0.002$	CTGTCACCTCT TTGTCACCTCT	$0.963 \\ 0.018$	$0.955 \\ 0.038$	$0.839 \\ 0.058$	$0.975 \\ 0.018$	$0.991 \\ 0.005$	$\begin{array}{c} 0.886 \\ 0.1111 \end{array}$
Position 242 (0.98)	EuA u 99	$_{00-04}^{\rm EuA}$	$_{05-09}^{\rm EuA}$	$^{\rm As}_{ m u \ 99}$	$_{00-04}^{\rm As}$	$^{\mathrm{As}}_{05-09}$	Position 293 (0.98)	EuA u 99	EuA 00-04	$_{05-09}^{\rm EuA}$	As u 99	$^{ m As}_{ m 00-04}$	$^{ m As}_{05-09}$
ATTT ATGTT	$0.945 \\ 0.002$	$0.949 \\ 0.005$	$0.779 \\ 0.13$	$0.984 \\ 0$	0.98 0	$0.988 \\ 0$	CGTAGACGGTTTGT CGTAGACGATTTGT CGTAGACGCTTTGT	$\begin{array}{c} 0.003 \\ 0.112 \\ 0.864 \end{array}$	$\begin{array}{c} 0.088 \\ 0.164 \\ 0.728 \end{array}$	$\begin{array}{c} 0.06 \\ 0.047 \\ 0.86 \end{array}$	$\begin{array}{c} 0.0033 \\ 0.022 \\ 0.952 \end{array}$	$\begin{array}{c} 0.005 \\ 0.002 \\ 0.951 \end{array}$	$\begin{array}{c} 0\\ 0\\ 0.978 \end{array}$
Position $323$ (0.95)	EuA u 99	$_{00-04}^{\rm EuA}$	$_{05-09}^{\rm EuA}$	$^{\rm As}_{ m u \ 99}$	$^{ m As}_{00-04}$	$^{\mathrm{As}}_{05-09}$	Position 354 (0.95)	EuA u 99	EuA 00-04	$_{05-09}^{\rm EuA}$	As u 99	$^{ m As}_{ m 00-04}$	$^{ m As}_{05-09}$
GGAAATGG GGGAACGG GGGAATGG	$\begin{array}{c} 0.06 \\ 0.058 \\ 0.846 \end{array}$	$\begin{array}{c} 0.346 \\ 0.01 \\ 0.621 \end{array}$	$\begin{array}{c} 0.578 \\ 0.002 \\ 0.418 \end{array}$	$\begin{array}{c} 0.015 \\ 0.073 \\ 0.871 \end{array}$	$\begin{array}{c} 0.04 \\ 0.002 \\ 0.938 \end{array}$	0.006 0.002 0.99	CAGT CGGT AAACT	$\begin{array}{c} 0.961 \\ 0.015 \\ 0.916 \end{array}$	$\begin{array}{c} 0.985 \\ 0.003 \\ 0.627 \end{array}$	$\begin{array}{c} 0.983 \\ 0.008 \\ 0.357 \end{array}$	$\begin{array}{c} 0.776 \\ 0.215 \\ 0.932 \end{array}$	$\begin{array}{c} 0.754 \\ 0.24 \\ 0.986 \end{array}$	$\begin{array}{c} 0.851 \\ 0.149 \\ 0.997 \end{array}$
Position 359 (0.98)	EuA u 99	${ m EuA}_{00-04}$	$_{05-09}^{\rm EuA}$	$^{\rm As}_{ m u \ 99}$	${}^{ m As}_{00-04}$	$^{\mathrm{As}}_{05-09}$	Position 368 (0.98)	EuA u 99	EuA 00-04	EuA $05-09$	As u 99	$^{ m As}_{ m 00-04}$	$^{ m As}_{05-09}$
AAATT	$0.024 \\ 0.057$	$0.113 \\ 0.256$	$0.061 \\ 0.573$	$0.014 \\ 0.053$	$0.011 \\ 0.003$	$0.001 \\ 0.002$	G	$0.974 \\ 0.026$	$0.973 \\ 0.027$	$0.937 \\ 0.062$	$0.997 \\ 0.002$	$0.903 \\ 0.097$	$0.671 \\ 0.329$

K. Bartoszek, P. Liò, A. Sorathiya

As 05-09 0.629 0.368

As 00-04 0.875 0.122

EuA 05-09 0.782 0.216

EuA 00-04 0.74 0.261

Position 389 (0.95)

As 05-09 0.111 0.889

As 00-04 0.003 0.997

As u 99 0.007 0.993

EuA 05–09 0.717 0.283

EuA 00-04 0.559 0.441

EuA u 99 0.206 0.791

٩Ŭ

Position 369(0.97)

As u 99 0.994 0.005

EuA u 99 0.969 0.023

> ACAT ACGT

432

102

TABLE VIII

D							i i	1						
	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.779 \\ 0.111 \\ 0.103 \end{array}$	$^{\mathrm{As}}_{05-09}$	0.985 0.015	$\mathbf{As}_{05-09}$	$0.984 \\ 0.001$	$^{\mathrm{As}}_{05-09}$	0.399 0.601 0	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.1111 \\ 0.887 \end{array}$	$\mathbf{As}_{05-09}$	$\begin{array}{c} 0.1111\\ 0\\ 0\\ 0\\ 0.868\\ 0.868\end{array}$	$^{ m As}_{05-09}$	$0.87 \\ 0.128$
	$\stackrel{ m As}{_{00-04}}$	$\begin{array}{c} 0.946 \\ 0.003 \\ 0.003 \end{array}$	$\mathop{\rm As}_{00-04}$	$0.971 \\ 0.029$	$\stackrel{ m As}{_{00-04}}$	$0.968 \\ 0.012$	${}^{ m As}_{ m 00-04}$	$\begin{array}{c} 0.106\\ 0.894\\ 0\end{array}$	$\stackrel{ m As}{_{00-04}}$	$0 \\ 0.998$	$\stackrel{ m As}{_{00-04}}$	$\begin{array}{c} 0 \\ 0 \\ 0 \\ 0 \\ 0.932 \end{array}$	$\stackrel{ m As}{_{00-04}}$	$0.948 \\ 0.043$
	As u 99	$\begin{array}{c} 0.978 \\ 0.001 \\ 0 \end{array}$	As u 99	$0.997 \\ 0.002$	As u 99	$0.941 \\ 0.016$	As u 99	$\begin{array}{c} 0.005 \\ 0.987 \\ 0.001 \end{array}$	As u 99	0.006 0.987	As u 99	$\begin{array}{c} 0.005 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0.95 \end{array}$	As u 99	$0.975 \\ 0.019$
	EuA $05-09$	$\begin{array}{c} 0.922 \\ 0.065 \\ 0 \end{array}$	EuA 05-09	$0.931 \\ 0.069$	EuA $05-09$	$0.404 \\ 0.518$	$E_{uA}$ 05-09	$\begin{array}{c} 0.062 \\ 0.369 \\ 0.538 \end{array}$	EuA $05-09$	$0.714 \\ 0.285$	EuA 05-09	$\begin{array}{c} 0.197 \\ 0.353 \\ 0.04 \\ 0.144 \\ 0.225 \end{array}$	EuA $05-09$	$0.961 \\ 0.037$
	EuA 00-04	$\begin{array}{c} 0.925 \\ 0.048 \\ 0 \end{array}$	$_{00-04}^{\rm EuA}$	$0.782 \\ 0.213$	EuA 00-04	0.735 0.203	EuA 00-04	$\begin{array}{c} 0.025 \\ 0.698 \\ 0.244 \end{array}$	EuA 00-04	$0.56 \\ 0.434$	EuA 00-04	$\begin{array}{c} 0.292 \\ 0.206 \\ 0.066 \\ 0.005 \\ 0.288 \end{array}$	EuA 00-04	$0.992 \\ 0.008$
	EuA u 99	$\begin{array}{c} 0.946 \\ 0.01 \\ 0.002 \end{array}$	EuA u 99	$0.904 \\ 0.096$	EuA u 99	$0.935 \\ 0.011$	EuA u 99	$\begin{array}{c} 0.029 \\ 0.945 \\ 0.01 \end{array}$	EuA u 99	$0.162 \\ 0.826$	EuA u 99	$\begin{array}{c} 0.279 \\ 0.013 \\ 0.033 \\ 0 \\ 0 \\ 0.573 \end{array}$	EuA u 99	$0.992 \\ 0.005$
	$\begin{array}{c} \text{Position 416} \\ (0.97) \end{array}$	CTCAG CTAAG CTCGG	Position 469 (0.98)	AG CG	Position 484 (1)	ACCAC ACTAC	Position 529 (0.96)	ATTGCCGA ATTGCTGA ATTGCAGA	Position 563 (0.99)	TGGC TGGT	Position 569 (0.96)	CTACCACCAA CTATCACCAA CTATCACCAA CTACTACCAA CCATCACCAA CAACAACCAA	Position $643$ (0.95)	ATGGC GTGGC
-	$^{\mathrm{As}}_{05-09}$	0.003 0.995	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0 \\ 0 \\ 0.004 \\ 0.98 \end{array}$	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.115 \\ 0.002 \\ 0.882 \end{array}$	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.729\\ 0\\ 0.112\\ 0.003\\ 0.067\end{array}$	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.111\\ 0.871\\ 0\end{array}$	$\overset{\mathrm{As}}{_{05-09}}$	0.864 0.119	$^{ m As}_{05-09}$	$0.873 \\ 0.119$
	${\mathop{\rm As}}_{00-04}$	$0.018 \\ 0.966$	${\mathop{\rm As}}_{00-04}$	$\begin{array}{c} 0 \\ 0.019 \\ 0.005 \\ 0.951 \end{array}$	${\mathop{\rm As}}_{00-04}$	$\begin{array}{c} 0.003 \\ 0 \\ 0.983 \end{array}$	${}^{ m As}_{ m 00-04}$	$\begin{array}{c} 0.978 \\ 0 \\ 0.006 \\ 0.003 \end{array}$	${\mathop{\rm As}}_{00-04}$	$\begin{array}{c} 0\\0.977\\0.001\end{array}$	${}^{ m As}_{ m 00-04}$	0.818 0.008	${}^{ m As}_{ m 00-04}$	$0.96 \\ 0.008$
	As u 99	0.097 0.880	As u 99	$\begin{array}{c} 0.002 \\ 0.0215 \\ 0.006 \\ 0.861 \end{array}$	As u 99	$\begin{array}{c} 0.009 \\ 0.006 \\ 0.981 \end{array}$	As u 99	$\begin{array}{c} 0.678 \\ 0 \\ 0.024 \\ 0.214 \\ 0.039 \end{array}$	As u 99	$\begin{array}{c} 0.031 \\ 0.952 \\ 0.001 \end{array}$	As u 99	$0.954 \\ 0.016$	As u 99	$0.937 \\ 0.016$
	EuA 05–09	$0.004 \\ 0.960$	EuA 05–09	$\begin{array}{c} 0.126\\ 0.417\\ 0.098\\ 0.33\end{array}$	EuA 05–09	$\begin{array}{c} 0.31 \\ 0.46 \\ 0.231 \end{array}$	EuA $05-09$	$\begin{array}{c} 0.215\\ 0.101\\ 0.649\\ 0.021\\ 0.004\end{array}$	EuA 05–09	$\begin{array}{c} 0.68 \\ 0.226 \\ 0.07 \end{array}$	EuA 05-09	$0.989 \\ 0.051$	EuA $05-09$	$0.941 \\ 0.051$
	EuA 00-04	$0.004 \\ 0.935$	EuA 00-04	$\begin{array}{c} 0.002 \\ 0.346 \\ 0.236 \\ 0.343 \end{array}$	EuA 00-04	$\begin{array}{c} 0.263 \\ 0.444 \\ 0.29 \end{array}$	EuA 00-04	$\begin{array}{c} 0.289\\ 0.241\\ 0.434\\ 0.005\\ 0.001\end{array}$	EuA 00-04	$\begin{array}{c} 0.557 \\ 0.275 \\ 0.158 \end{array}$	EuA 00-04	0.999 0.003	${ m EuA}_{00-04}$	$0.993 \\ 0.003$
	EuA u 99	$0.021 \\ 0.937$	EuA u 99	$\begin{array}{c} 0 \\ 0.192 \\ 0.084 \\ 0.651 \end{array}$	EuA u 99	$\begin{array}{c} 0.157 \\ 0.237 \\ 0.601 \end{array}$	EuA u 99	$\begin{array}{c} 0.146\\ 0.086\\ 0.294\\ 0.281\\ 0.002\end{array}$	EuA u 99	$\begin{array}{c} 0.369 \\ 0.583 \\ 0.044 \end{array}$	EuA u 99	0.987 0.002	EuA u 99	$0.982 \\ 0.002$
	$\begin{array}{c} \text{Position 392} \\ (0.96) \end{array}$	TTCCACGG TTCCATGG	Position 429 (0.97)	CCGGTGCACTCGC CCGGTGCACTTGC CTGGTGCGCTTGC CTGGTGCGCTTGC CTGGTGCACTTGC	Position 479 (0.97)	CAGT CGGT CTGT	Position 496 (1)	GCATTTGG GCTCTTGG GCTTTTGG GCTTTTGG GCCTTTGG GCGTTTGG	Position 538 (0.97)	TCACA TCCCA GCCCA	Position 592 (1)	CATGA CACGA	Position 625 (0.96)	GCTAA GCAAA

# Influenza Differentiation and Evolution

$^{\mathrm{As}}_{05-09}$	$0.189 \\ 0.811$	$\stackrel{\mathrm{As}}{_{05-09}}$	0.99 0.005	$^{\mathrm{As}}_{05-09}$	0.966 0	$\stackrel{ m As}{_{05-09}}$	$\begin{array}{c} 0\\ 0.001\\ 0.998\end{array}$	$\stackrel{ m As}{_{05-09}}$	$\begin{array}{c} 0.121\\0\\0.875\end{array}$	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.121\\0\\0.875\end{array}$	$^{\mathrm{As}}_{05-09}$	$0.998 \\ 0.002$	$\overset{\mathrm{As}}{_{05-09}}$	$0.001 \\ 0.982$
$\mathbf{A_S}_{00-04}$	.003 $0.995$	${}^{ m As}_{ m 00-04}$	$0.983 \\ 0$	${}^{ m As}_{ m 00-04}$	$0.995 \\ 0$	$\operatorname*{As}_{00-04}$	$\begin{array}{c} 0.005 \\ 0.012 \\ 0.983 \end{array}$	${}^{ m As}_{ m 00-04}$	$\begin{array}{c} 0.018 \\ 0 \\ 0.98 \end{array}$	${}^{ m As}_{ m 00-04}$	$\begin{array}{c} 0.018 \\ 0 \\ 0.98 \end{array}$	${}^{ m As}_{ m 00-04}$	$0.997 \\ 0$	${}^{ m As}_{ m 00-04}$	$0.002 \\ 0.974$
As u 99	$\begin{array}{c} 0.01 \ 0 \\ 0.988 \end{array}$	As u 99	$0.923 \\ 0.016$	As u 99	$0.964 \\ 0.033$	As u 99	$\begin{array}{c} 0.014 \\ 0.014 \\ 0.972 \end{array}$	As u 99	$\begin{array}{c} 0\\ 0.007\\ 0.991 \end{array}$	As u 99	$\begin{array}{c} 0\\ 0.007\\ 0.991 \end{array}$	As u 99	$0.955 \\ 0.003$	As u 99	0.085 0.890
EuA 05-09	$0.591 \\ 0.409$	$_{05-09}^{\rm EuA}$	$0.448 \\ 0.528$	$_{05-09}^{\rm EuA}$	$0.448 \\ 0.52$	$_{05-09}^{\rm EuA}$	$\begin{array}{c} 0.094 \\ 0.026 \\ 0.877 \end{array}$	$_{05-09}^{\rm EuA}$	$\begin{array}{c} 0.033 \\ 0.0427 \\ 0.924 \end{array}$	$_{05-09}^{\rm EuA}$	$\begin{array}{c} 0.033 \\ 0.0427 \\ 0.924 \end{array}$	$_{05-09}^{\rm EuA}$	$0.896 \\ 0.103$	$_{05-09}^{\rm EuA}$	0.013 0.945
EuA 00-04	$0.242 \\ 0.758$	EuA 00-04	$0.745 \\ 0.227$	EuA 00-04	$0.734 \\ 0.245$	$_{ m EuA}_{ m 00-04}$	$\begin{array}{c} 0.127 \\ 0.079 \\ 0.787 \end{array}$	EuA 00-04	$\begin{array}{c} 0.031 \\ 0.0544 \\ 0.911 \end{array}$	EuA 00-04	$\begin{array}{c} 0.031 \\ 0.0544 \\ 0.911 \end{array}$	EuA 00-04	$0.902 \\ 0.098$	EuA 00-04	$0.008 \\ 0.964$
EuA u 99	$0.008 \\ 0.992$	EuA u 99	$0.154 \\ 0.016$	EuA u 99	$0.984 \\ 0.01$	EuA u 99	$\begin{array}{c} 0.045 \\ 0.018 \\ 0.930 \end{array}$	EuA u 99	$\begin{array}{c} 0.005 \\ 0.05 \\ 0.937 \end{array}$	EuA u 99	$\begin{array}{c} 0.005 \\ 0.05 \\ 0.937 \end{array}$	EuA u 99	$0.987 \\ 0.002$	EuA u 99	$0.023 \\ 0.948$
Position 689 (0.97)	A G	Position 736 (0.98)	CCTAG CCTAA	Position 761 (0.99)	ATGA ATAA	Position 785 (0.99)	CTTA CATA CCTA	Position 811 (0.96)	CTGCA ATACA ATGCA	Position 811 (0.96)	CTGCA ATACA ATGCA	Position 833 (1)	CT AT	Position $840$ (0.97)	CTGC TTGC
$^{\mathrm{As}}_{05-09}$	$0.111 \\ 0.002$	$^{\mathrm{As}}_{05-09}$	$0.111 \\ 0.889$	$^{\mathrm{As}}_{05-09}$	0.003 0.994	$^{\mathrm{As}}_{05-09}$	$0.985 \\ 0$	$^{\mathrm{As}}_{05-09}$	$0.127 \\ 0.873$	$^{\mathrm{As}}_{05-09}$	$0.127 \\ 0.873$	$^{\mathrm{As}}_{05-09}$	$0.072 \\ 0.893$	$^{\mathrm{As}}_{05-09}$	0.232 0.759
$\mathbf{As}_{00-04}$	$0 \\ 0.006$	$\mathop{\mathrm{As}}\limits_{00-04}$	0.003 0.992	$\stackrel{ m As}{_{00-04}}$	$\begin{array}{c} 0.017 \\ 0.974 \end{array}$	$\mathop{\mathrm{As}}_{00-04}$	0.997 0	$\mathop{\mathrm{As}}_{00-04}$	0.018 0.982	${\mathop{\rm As}}_{00-04}$	$0.018 \\ 0.982$	${\mathop{\rm As}}_{00-04}$	$0.132 \\ 0.866$	$\mathop{\mathrm{As}}\limits_{00-04}$	$0.034 \\ 0.963$
As u 99	0.006 0.006	As u 99	$0.014 \\ 0.924$	As u 99	$\begin{array}{c} 0.062 \\ 0.928 \end{array}$	As u 99	$0.961 \\ 0.001$	As u 99	$0.002 \\ 0.999$	As u 99	0.002 0.999	As u 99	$\begin{array}{c} 0.001 \\ 0.987 \end{array}$	As u 99	$0.004 \\ 0.987$
EuA 05–09	0.71415 0.55081	$_{05-09}^{\rm EuA}$	$0.066 \\ 0.929$	$_{05-09}^{\rm EuA}$	$0.01 \\ 0.967$	${ m EuA}_{05-09}$	$0.859 \\ 0.109$	$_{05-09}^{\rm EuA}$	$0.002 \\ 0.998$	$_{05-09}^{\rm EuA}$	$0.002 \\ 0.998$	$_{05-09}^{\rm EuA}$	$0.008 \\ 0.991$	$_{05-09}^{\rm EuA}$	0.095 0.830
EuA 00-04	$0.55962 \\ 0.35042$	EuA 00-04	$0.051 \\ 0.938$	EuA 00-04	0.023 0.923	$_{00-04}^{\rm EuA}$	$0.742 \\ 0.24$	EuA 00-04	$0.01 \\ 0.990$	EuA 00-04	$0.01 \\ 0.990$	EuA 00-04	0.03 0.949	EuA 00-04	$0.04 \\ 0.938$
EuA u 99	$0.162 \\ 0.056$	EuA u 99	0.073 0.903	EuA u 99	$0.044 \\ 0.948$	EuA u 99	0.886 0.09	EuA u 99	$0.002 \\ 0.995$	EuA u 99	$0.002 \\ 0.995$	EuA u 99	$0.016 \\ 0.977$	EuA u 99	$0.021 \\ 0.159$
Position 683 (0.98)	TTGC TCGC	Position 709 (0.96)	CATGC CAGGC	Position 751 (0.96)	GGTTT GGTCT	Position 769 (1)	CTTGAAAATTT ATTGAAAATTT	Position 796 (0.98)	AG CG	Position 796 (0.98)	AG CG	$\begin{array}{c} \text{Position 820} \\ (0.97) \end{array}$	TTCAAA TTCAAG	Position 837 (0.95)	TT

$^{\rm As}_{4 \ 05-09}$	$\begin{array}{ccc} 2 & 0.001 \\ 4 & 0.982 \end{array}$	4 05-09	8 0.889 0.11	4 05–09	$\begin{array}{ccc} 2 & 0.744 \\ 8 & 0.255 \end{array}$	4 05–09	$\begin{array}{ccc} 2 & 0.402 \\ 0.328 \\ 5 & 0.256 \end{array}$	4 05-09	2 0.868 2 0.11868	4 05–09	8 0.994 6 0.003 2 0.001	4 05–09	0 7 0.998		
$^{\mathrm{As}}_{00-0}$	0.00.0	As 00-00	0.99 0	$^{\mathrm{As}}_{00-0}$	0.85 0.14	$_{00-0}^{\mathrm{As}}$	$\begin{array}{c} 0.77\\ 0.1\\ 0.1\\ 0.11.\end{array}$	As 00-0	0.02	$_{00-0}^{\mathrm{As}}$	0.94 0.04 0.00	$_{00-0}^{\mathrm{As}}$	0.97		
As u 99	$\begin{array}{c} 0.085 \\ 0.890 \end{array}$	As u 99	$0.993 \\ 0.001$	As u 99	$0.968 \\ 0.027$	$_{ m u}^{ m As}$ 99	$\begin{array}{c} 0.986\\ 0.007\\ 0\end{array}$	As u 99	0.953 0.02	As u 99	$\begin{array}{c} 0.899\\ 0.025\\ 0.011\end{array}$	As u 99	$\begin{array}{c} 0.008 \\ 0.942 \end{array}$		
EuA 05-09	$\begin{array}{c} 0.013 \\ 0.945 \end{array}$	EuA 05–09	$0.956 \\ 0.038$	EuA 05-09	$\begin{array}{c} 0.842 \\ 0.151 \end{array}$	EuA $05-09$	$\begin{array}{c} 0.774 \\ 0.062 \\ 0.149 \end{array}$	EuA 05-09	0.972 0	EuA 05–09	$\begin{array}{c} 0.362 \\ 0.279 \\ 0.355 \end{array}$	EuA 05-09	$0.01 \\ 0.953$		
EuA 00-04	$\begin{array}{c} 0.008 \\ 0.964 \end{array}$	EuA 00-04	$^{0.99}_{0}$	EuA 00-04	$\begin{array}{c} 0.802 \\ 0.184 \end{array}$	EuA 00-04	$\begin{array}{c} 0.797 \\ 0.025 \\ 0.177 \end{array}$	EuA 00-04	$0.983 \\ 0.001$	EuA 00-04	$\begin{array}{c} 0.445 \\ 0.371 \\ 0.177 \end{array}$	EuA 00-04	$\begin{array}{c} 0.084 \\ 0.889 \end{array}$		
EuA u 99	$\begin{array}{c} 0.023 \\ 0.948 \end{array}$	EuA u 99	$0.971 \\ 0$	EuA u 99	$0.987 \\ 0.002$	EuA u 99	$\begin{array}{c} 0.104 \\ 0.028 \\ 0.052 \end{array}$	EuA u 99	$0.974 \\ 0.003$	EuA u 99	$\begin{array}{c} 0.787 \\ 0.195 \\ 0.006 \end{array}$	EuA u 99	$\begin{array}{c} 0.026 \\ 0.938 \end{array}$		
Position 840 (0.97)	CTGC TTGC	Position 885 (1)	TTGATCG CTGATCG	$\begin{array}{c} \text{Position 908} \\ (0.97) \end{array}$	AT GT	Position 933 (0.96)	TGAAA TAAAA TTAAA	Position 965 (0.97)	GAGTC AAGTC	$\begin{array}{c} \text{Position 980} \\ (0.97) \end{array}$	GAATATC GAGTATC GAGTACC	Position 993 (1)	AGCAACAG AACAGCAG		
$\stackrel{\mathrm{As}}{_{05-09}}$	$0.232 \\ 0.759$	$^{\mathrm{As}}_{05-09}$	$0.328 \\ 0.67$	$^{\mathrm{As}}_{05-09}$	$0.671 \\ 0.325$	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.038 \\ 0.006 \\ 0.954 \end{array}$	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0 \\ 0 \\ 0.239 \\ 0.685 \end{array}$	$^{\mathrm{As}}_{05-09}$	$0.129 \\ 0.867$	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.014 \\ 0.984 \end{array}$	$^{\mathrm{As}}_{05-09}$	$0.692 \\ 0.297$
${\mathop{\rm As}}_{00-04}$	$0.034 \\ 0.963$	$^{\mathrm{As}}_{00-04}$	$\begin{array}{c} 0.092 \\ 0.903 \end{array}$	$\mathop{\rm As}_{00-04}$	$0.902 \\ 0.091$	${\mathop{\rm As}}_{00-04}$	$\begin{array}{c} 0.134 \\ 0 \\ 0.862 \end{array}$	$\stackrel{ m As}{_{00-04}}$	$\begin{array}{c} 0 \\ 0 \\ 0.123 \\ 0.802 \end{array}$	$\stackrel{ m As}{_{00-04}}$	$0.054 \\ 0.945$	$\stackrel{ m As}{_{00-04}}$	$0.017 \\ 0.969$	$\stackrel{\mathrm{As}}{_{00-04}}$	$0.995 \\ 0.0015$
As u 99	$\begin{array}{c} 0.004 \\ 0.987 \end{array}$	As u 99	$0 \\ 0.966$	As u 99	$0.993 \\ 0$	As u 99	$\begin{array}{c} 0.006 \\ 0.053 \\ 0.928 \end{array}$	As u 99	$\begin{array}{c} 0 \\ 0 \\ 0.002 \\ 0.929 \end{array}$	As u 99	$\begin{array}{c} 0.022 \\ 0.978 \end{array}$	As u 99	$\begin{array}{c} 0.092 \\ 0.903 \end{array}$	As u 99	$0.99 \\ 0.001$
EuA 05–09	$\begin{array}{c} 0.095 \\ 0.830 \end{array}$	EuA 05-09	$\begin{array}{c} 0.062 \\ 0.895 \end{array}$	EuA 05–09	$0.907 \\ 0.057$	EuA 05–09	$\begin{array}{c} 0.001 \\ 0.002 \\ 0.997 \end{array}$	EuA 05–09	$\begin{array}{c} 0.07 \\ 0.084 \\ 0.138 \\ 0.543 \end{array}$	EuA 05–09	$\begin{array}{c} 0.004 \\ 0.994 \end{array}$	EuA 05–09	$\begin{array}{c} 0.01 \\ 0.985 \end{array}$	EuA 05–09	$0.972 \\ 0.021$
EuA 00-04	$0.04 \\ 0.938$	EuA 00-04	$\begin{array}{c} 0.021 \\ 0.934 \end{array}$	EuA 00-04	$0.973 \\ 0.024$	EuA 00-04	$\begin{array}{c} 0.019 \\ 0.001 \\ 0.978 \end{array}$	EuA 00-04	$\begin{array}{c} 0.097 \\ 0.061 \\ 0.135 \\ 0.599 \end{array}$	EuA 00-04	$\begin{array}{c} 0.04 \\ 0.957 \end{array}$	EuA 00-04	$\begin{array}{c} 0.001 \\ 0.994 \end{array}$	EuA 00-04	$0.992 \\ 0.003$
EuA u 99	$\begin{array}{c} 0.021 \\ 0.159 \end{array}$	EuA u 99	$\begin{array}{c} 0.028 \\ 0.942 \end{array}$	EuA u 99	$\begin{array}{c} 0.961 \\ 0.028 \end{array}$	EuA u 99	$\begin{array}{c} 0.01 \\ 0.119 \\ 0.867 \end{array}$	EuA u 99	$\begin{array}{c} 0 \\ 0 \\ 0 \\ 0.930 \end{array}$	EuA u 99	$\begin{array}{c} 0.01 \\ 0.977 \end{array}$	EuA u 99	$\begin{array}{c} 0.013 \\ 0.974 \end{array}$	EuA u 99	0.997
Position 837 (0.95)	TC	Position 852 (0.95)	TAAT TCAT	Position 899 (0.99)	TTCAAAT TCCAAAA	Position 923 (0.98)	AGA GAA AAA	Position 935 (0.95)	AAAAGAGGGCCTGCTACGGCAGG AAAAGAGGGGCCTTCTACGGCAGG AAAAGAGGCCCTTCTACGGGAGG AAAAGAGGCCCTTCTACGGAAGG	Position 971 (0.96)	ATGAGGGAGGA ATGAGGGAAGA	Position 990 (0.98)	AGAA AGGA	Position 1023 (1)	GTCAT ATCAT

# 4.4. NA results

The NA gene segment sequences were analyzed to find differences between hosts, regions, time periods and geographical and temporal interactions. The results are presented in Fig. 5 and Tables IX, X, XI and XIII.



Fig. 5. Top left: scores along the alignment of the NA gene segment for splits between hosts, cut-off for graph is 0.95. Top scoring examples are in Table XII. Top right: scores along the alignment of the NA gene segment for splits between regions, cut-off for graph is 0.95. Top scoring examples are in Table X. Bottom left: scores along the alignment of the NA gene segment for splits between time periods, cut-off for graph is 0.95. Top scoring examples are in Table XI. Bottom right: scores along the alignment of the NA gene segment for splits between time periods, cut-off for graph is 0.95. Top scoring examples are in Table XI. Bottom right: scores along the alignment of the NA gene segment for splits between regions combined with time periods, cut-off for graph is 0.95. Top scoring examples are in Table XI. Bottom XI. Bottom right: scores along the alignment of the NA gene segment for splits between regions combined with time periods, cut-off for graph is 0.95. Top scoring examples are in Table XI.

TABLE IX

Top scoring (above 0.85) region specific changes of NA gene segment. Regions with gaps were included. Even though regions with gaps are not visible, they were present in sequences that had too few representatives, excluding them would leave only positions 328, 504, 512, 1191, 1249, 1300.

Position 94 Score 0.98	Eurasia	Americas	Position 201 Score 0.96	Eurasia	Americas
T C	$\begin{array}{c} 0.9 \\ 0.088 \end{array}$	$0.995 \\ 0.001$	Ā	$0.299 \\ 0.675$	$0.002 \\ 0.983$
Position 328 Score 0.91	Eurasia	Americas	Position 504 Score 0.96	Eurasia	Americas
T C	$0.009 \\ 0.983$	$0.087 \\ 0.913$	T C	$0.857 \\ 0.141$	$\begin{array}{c} 0.995\\ 0.004\end{array}$
Position 512 Score 0.92	Eurasia	Americas	Position 1191 Score 0.85	Eurasia	Americas
AG CG	0.991 0.009	$0.894 \\ 0.106$	AG TC AA GA	$\begin{array}{c} 0.81 \\ 0.025 \\ 0.095 \\ 0.062 \end{array}$	$0.951 \\ 0.034 \\ 0.009 \\ 0$
Position 1249 Score 0.89	Eurasia	Americas	Position 1300 Score 0.96	Eurasia	Americas
A C	$0.013 \\ 0.984$	$\begin{array}{c} 0.1 \\ 0.9 \end{array}$	GGTC CGTC	$0.903 \\ 0.064$	$\begin{array}{c} 0.975 \\ 0 \end{array}$
Position 1493 Score 0.94	Eurasia	Americas	Position 1497 Score 0.83	Eurasia	Americas
GG	0.086 0.901	$0.007 \\ 0.99$	C	0.09 0.907	$0.008 \\ 0.992$

Top scoring (above 0.85) time period spe 340, 507, 698, 548 and 1184 would remai Position 201	cific ch n. Until 3	anges 2000–	of NA 2005-	e gene segmen Position 340	t. If re Until	gions 2000-	with g <sup>2005-</sup>	aps were excl Position 507	Until	2000–2000–	2005-
Score: 0.9	1999	2004	2009	Score: 0.92	1999	2004	2009	Score: 0.99	1999	2004	2009
¥	0.005 0.959	0.116 0.86	0.258 0.733	A	$0.848 \\ 0.117$	0.951 0.03	0.967	A	$0.998 \\ 0.001$	$0.999 \\ 0.001$	0.882 0.117
Position 548 Score: 0.92	Until 3 1999	2000- 2004	2005 - 2009	Position 698 Score: 0.88	Until 1999	2000 - 2004	2005 - 2009				
AG GG	0.005 0.974	$0.104 \\ 0.891$	0.233 0.76	AC GA GG	$\begin{array}{c} 0.073 \\ 0.925 \\ 0.001 \end{array}$	$\begin{array}{c} 0.026 \\ 0.97 \\ 0.001 \end{array}$	$\begin{array}{c} 0.031 \\ 0.859 \\ 0.107 \end{array}$				
Position 1184 Score: 0.85	Until : 1999	2000- 2004	2005 - 2009	Position 1558 Score: 0.88	Until 1999	2000 - 2004	2005 - 2009				
AACTAAAAGCATTAGTTCAAGAAAC AACCAAAAGCACTAATTCCAGGAGC AACGATCAGCGAGAAGTTACGCTCA AACTAAAAGTAACAGAGTTAGGCTCA GACTAAAAGTAACAGACTTAGAAAG GACTAAAAGTAACAGAGTTAGAAAG AACGATCAGCGAGAAGTCACGCTTA GACCAAAAGTAACAGACTTAGAAAG	$\begin{array}{c} 0 \\ 0.003 \\ 0.101 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array}$	$\begin{array}{c} 0 \\ 0.107 \\ 0.389 \\ 0.001 \\ 0.082 \\ 0 \\ 0 \\ 0 \\ 0.001 \end{array}$	$\begin{array}{c} 0.06\\ 0.159\\ 0.097\\ 0.068\\ 0.016\\ 0.1\\ 0.1\\ 0.094 \end{array}$	TTCTACT TCTACT- TCTAC- TCTAC-	$\begin{array}{c} 0.074 \\ 0.039 \\ 0.845 \\ 0.025 \end{array}$	$\begin{array}{c} 0.052\\ 0.0432\\ 0.850\\ 0.04\end{array}$	$\begin{array}{c} 0.01\\ 0.008\\ 0.961\\ 0.009\end{array}$				

TABLE X

438

# K. Bartoszek, P. Liò, A. Sorathiya

IC	$^{84}$	
We	10	
aces	990,	
pl	ين. ت	
uch	. 8	
08 s	755	
all 2	393,	
н.	33, (	
all	ы м	
lent,	328	
egm	ted:	
ne s	esen	
ge.	pre	
NA	hese	
s of	of t	
nge	'n,	
$_{\rm cha}$	ema	
iffic	ld r	
spec	mom	
ost s	' SU	
) he	sitic	
0.85	bo	
ove	127	
(abc	$\operatorname{red}$	
ng	gno	
cori	are i	384.
op s	bs 8	d 1:
le t	1 ga	) an
of tł	witl	078
ng (	$\operatorname{suc}$	of 1
esti	regi	art
nter	If 1	dqn
ist i	nd.	as
$M_{C}$	fou	(as)

TABLE XI

Position 96 Score 0.86	Swine	Avian	Equine	Human Score 0.96	Position 151	Swine	Avian	Equine	Human
GGAATGGC-TAACTTAATATTACAA	0.002	0	0	0.045	Α	0.145	0.273	0.101	0.007
GGGATAAT-TAGTCTAATGTTGCAA	0.002	0	0	0.04	L	0.009	0.113	0.876	0
TCCACAAT-ATGCTTCTTCATGCAA	0.004	0	0	0.217	IJ	0.844	0.597	0	0.992
GCAACAGT-ATGTTTCCTCATGCAG	0.002	0.079	0	0					
GGAATAGT-TAGCTTAATGTTACAA	0.004	0.196	0	0.029					
GGAATAAT-TAGTCTAATGTTGCAA	0.002	0.001	0	0.15					
GCCACAAT-ATGCTTCCTTATGCAA	0.067	0.002	0	0.269					
Position 171	Swine	Avian	Equine	Human	Position 265	Swine	Avian	Equine	Human
Score 0.86				Score 0.88					
	0.004	0	0	0.23		0.002	0.058	0	0.008
	0.013	0.001	0	0.149	AACAC	0.388	0.02	0	0.943
	0.007	0.105	0	0.006	AATGT	0.002	0.074	0.011	0
	0	0.062	0	0.003	AGCAC	0.022	0.08	0	0
	0.002	0	0	0.046	AATAC	0.414	0.488	0.427	0.042
	0	0	0	0.068	GAACC	0	0.042	0	0
	0.052	0	0	0.24	AACAA	0.03	0.051	0	0
Position 328	Swine	Avian	Equine	Human	Position 386	Swine	Avian	Equine	Human
Score $0.96$				Score 0.99					
£	0.002	0.104	0.011	0	Α	0.013	0.162	0.101	0.001
C	0.998	0.884	0.978	1	С	0.987	0.838	0.899	0.999
Position 533	Swine	Avian	Equine	Human	Position 657	Swine	Avian	Equine	Human
Score 0.98				Score $0.94$					
LL	0.056	0.29	0	0.001	Α	0.937	0.698	0.798	0.989
AT	0.87	0.682	1	0.991	C	0.022	0.064	0.202	0
					U	0.041	0.236	0	0.012

# Influenza Differentiation and Evolution

Human	$\begin{array}{c} 0\\ 0.963\\ 0.036\end{array}$	Human	$\begin{array}{c} 0.001 \\ 0.003 \\ 0.956 \\ 0.011 \\ 0.03 \end{array}$	Human	0.017 0.983	Human	$\begin{array}{c} 0.964 \\ 0.001 \\ 0.035 \end{array}$
Equine	$\begin{array}{c} 0.809 \\ 0 \\ 0.112 \end{array}$	Equine	$\begin{array}{c} 0 \\ 0.1 \\ 0 \\ 0 \\ 0.9 \end{array}$	Equine	$0.112 \\ 0.888$	Equine	$\begin{array}{c} 0.101\\ 0.899\\ 0\end{array}$
Avian	$\begin{array}{c} 0.172\\ 0.65\\ 0.17\end{array}$	Avian	$\begin{array}{c} 0.051 \\ 0.21 \\ 0.301 \\ 0.078 \\ 0.272 \end{array}$	Avian	$0.207 \\ 0.792$	Avian	$\begin{array}{c} 0.557 \\ 0.088 \\ 0.325 \end{array}$
Swine	$\begin{array}{c} 0.022 \\ 0.94 \\ 0.035 \end{array}$	Swine	$\begin{array}{c} 0.002 \\ 0.095 \\ 0.85 \\ 0 \\ 0 \\ 0.043 \end{array}$	Swine	$0.082 \\ 0.915$	Swine	$\begin{array}{c} 0.881 \\ 0.013 \\ 0.106 \end{array}$
Position 815	A C	Position 818	AAC CAC TAC CAT TAT	Position 903	GAAGA GAGGA	Position 1008	A C Q
Human Score 0.86	10	Human Score 0.87	0.006 0.994	Human Score 0.88	$\begin{array}{c} 0.001 \\ 0.02 \\ 0.978 \\ 0 \end{array}$	Human Score 0.85	0.017 0.001 0.981 0.001
Equine	$0.1 \\ 0.9$	Equine	0.9 0.1	Equine	$\begin{array}{c} 0\\ 0.1\\ 0\\ 0\end{array}$	Equine	$\begin{array}{c} 0.101 \\ 0.876 \\ 0 \\ 0.022 \end{array}$
Avian Equine	$\begin{array}{ccc} 0.85 & 0.1 \\ 0.15 & 0.9 \end{array}$	Avian Equine	$\begin{array}{ccc} 0.517 & 0.9 \\ 0.482 & 0.1 \end{array}$	Avian Equine	$\begin{array}{cccc} 0.203 & 0\\ 0.248 & 0.1\\ 0.444 & 0\\ 0.06 & 0 \end{array}$	Avian Equine	$\begin{array}{cccc} 0.324 & 0.101 \\ 0.17 & 0.876 \\ 0.444 & 0 \\ 0.063 & 0.022 \end{array}$
Swine Avian Equine	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	Swine Avian Equine	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	Swine Avian Equine	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	Swine Avian Equine	$\begin{array}{cccccccccccccccccccccccccccccccccccc$

440

# K. Bartoszek, P. Liò, A. Sorathiya

Position 1078 Score 0.85	Swine	Avian	Equine	Human Score 0.88	Position 1117	Swine	Avian	Equine	Human
GCTTCAGCA—GTAGCCATTGCCT GAGAGGGCA—GCTGTAATCCAGT GCTCCAGCA—GTAGCCATTGTTT GCTCCAGCA—GTAGCCATTGTTT GCTCCAGCA—GTAGCCATTGCTT GCTCCAGCA—GTAGCCATTGCTT GAACGGGTA—GTTGTGGTCCGGGT GAAAGGGCA—GCTGTAATCCAGT GAGAGGGCA—GCTGTAATCCAGT AGACAGGCA—GTTGTGGTCCAGT	$\begin{array}{c} 0.024\\ 0\\ 0\\ 0.017\\ 0.041\\ 0.007\\ 0.002\\ 0\\ 0\\ 0\end{array}$	$\begin{array}{c} 0.002\\ 0.001\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0 \end{array}$	000000000	$\begin{array}{c} 0.108\\ 0.052\\ 0.072\\ 0.07\\ 0.234\\ 0.234\\ 0.019\\ 0.073\\ 0.0816\\ 0.046\end{array}$	AG GG GT	$\begin{array}{c} 0.002\\ 0.959\\ 0.026\\ 0.009\end{array}$	$\begin{array}{c} 0.07 \\ 0.362 \\ 0.475 \\ 0.073 \end{array}$	$\begin{array}{c} 0.112\\ 0\\ 0\\ 0\end{array}$	0 0.96 0.04 0
Position 1299 Score 0.93	Swine	Avian	Equine	Human Score 0.97	Position 1319	Swine	Avian	Equine	Human
A T	$0.408 \\ 0.577$	0.007 0.986	0	$0.6 \\ 0.381$	TTT CTT ATT	$\begin{array}{c} 0.883 \\ 0.115 \\ 0 \end{array}$	$\begin{array}{c} 0.701 \\ 0.207 \\ 0.06 \end{array}$	$\begin{array}{c} 1\\ 0\\ 0 \end{array}$	$\begin{array}{c} 0.996 \\ 0.004 \\ 0 \end{array}$
Position 1384 Score 0.94	Swine	Avian	Equine	Human Score 0.97	Position 1433	Swine	Avian	Equine	Human
AT GT	0.989 0.011	0.989 0.006	$\begin{array}{c} 1\\ 0 \end{array}$	$0.751 \\ 0.248$	CAG TAG CTC	$\begin{array}{c} 0.887 \\ 0.104 \\ 0.004 \end{array}$	$\begin{array}{c} 0.584 \\ 0.341 \\ 0.077 \end{array}$	$\begin{array}{c} 0.101\\ 0\\ 0.9\end{array}$	$\begin{array}{c} 0.993 \\ 0.004 \\ 0 \end{array}$

# Influenza Differentiation and Evolution

Η	
$\times$	
Ē	
닖	
7	
Ĥ	

442

Top scoring (above 0.85) time period and region specific changes of NA gene segment. If regions with gaps are ignored only 319, 328, 403, 408, 457, 497, 504, 507, 512, 548, 555, 560, 626, 677, 698, 701, 704, 761 and 788 would remain.

Position 9 Score: 0.96	EuA u 99	EuA 00-04	${ m EuA}_{05-09}$	As u 99	$\mathop{\mathrm{As}}\limits_{00-04}$	$^{ m As}_{05-09}$	Position 46 Score: 0.86	EuA u 99	${ m EuA}_{00-04}$	EuA $05-09$	As u 99	${}^{ m As}_{ m 00-04}$	$\stackrel{ m As}{_{05-09}}$
CAAAAGCA 	$\begin{array}{c} 0.73\\ 0.009\\ 0.041\\ 0.092\\ 0.016\end{array}$	$\begin{array}{c} 0.788 \\ 0.001 \\ 0.079 \\ 0.05 \\ 0.019 \end{array}$	$\begin{array}{c} 0.77 \\ 0 \\ 0.085 \\ 0.038 \\ 0.004 \end{array}$	$\begin{array}{c} 0.706 \\ 0.022 \\ 0.027 \\ 0.084 \\ 0.057 \end{array}$	$\begin{array}{c} 0.622\\ 0.148\\ 0.019\\ 0.058\\ 0.021\\ 0.021 \end{array}$	$\begin{array}{c} 0.932\\ 0.004\\ 0.011\\ 0.006\\ 0.017\\ \end{array}$	AATCCAAACCA AACCCAAATCA AATCCAAATCA	$\begin{array}{c} 0.061 \\ 0.014 \\ 0.009 \\ 0.879 \end{array}$	$\begin{array}{c} 0.049 \\ 0.042 \\ 0.025 \\ 0.774 \end{array}$	$\begin{array}{c} 0.024 \\ 0.035 \\ 0.057 \\ 0.788 \end{array}$	$\begin{array}{c} 0.003 \\ 0.007 \\ 0.004 \\ 0.875 \end{array}$	$\begin{array}{c} 0.018 \\ 0.013 \\ 0.006 \\ 0.869 \end{array}$	$\begin{array}{c} 0.052 \\ 0.107 \\ 0.127 \\ 0.644 \end{array}$
Position 62 Score: 0.85	EuA u 99	EuA 00-04	${ m EuA}_{05-09}$	As u 99	${}^{ m As}_{ m 00-04}$	$^{ m As}_{05-09}$	Position 72 Score: 0.85	EuA u 99	EuA 00-04	EuA $05-09$	As u 99	${}^{ m As}_{ m 00-04}$	$\stackrel{ m As}{_{05-09}}$
-TAATAGC -TAATAAC	$0.103 \\ 0.773$	$0.193 \\ 0.677$	0.087 0.853	$\begin{array}{c} 0.03 \\ 0.864 \end{array}$	0.007	$0.002 \\ 0.905$	TCGG TTGG TCTC	0.006 0.906 0.008	$\begin{array}{c} 0.2\\ 0.747\\ 0\end{array}$	$\begin{array}{c} 0.271 \\ 0.694 \\ 0.001 \end{array}$	$\begin{array}{c} 0.049 \\ 0.855 \\ 0.043 \end{array}$	$\begin{array}{c} 0.044 \\ 0.884 \\ 0.053 \end{array}$	$\begin{array}{c} 0.025 \\ 0.898 \\ 0.042 \end{array}$
Position 77 Score: 0.86	EuA u 99	EuA 00-04	$_{05-09}^{\rm EuA}$	As u 99	${}^{ m As}_{ m 00-04}$	$^{ m As}_{05-09}$	Position 94 Score: 0.99	EuA u 99	$_{00-04}^{\rm EuA}$	EuA $05-09$	As u 99	${}^{ m As}_{ m 00-04}$	${}^{ m As}_{05-09}$
09-19-1-	$\begin{array}{c} 0.032 \\ 0.036 \\ 0.887 \end{array}$	$\begin{array}{c} 0.02 \\ 0.033 \\ 0.927 \end{array}$	$\begin{array}{c} 0.019 \\ 0.007 \\ 0.954 \end{array}$	$\begin{array}{c} 0.067 \\ 0.038 \\ 0.778 \end{array}$	$\begin{array}{c} 0.032 \\ 0.051 \\ 0.825 \end{array}$	$\begin{array}{c} 0.047 \\ 0.018 \\ 0.906 \end{array}$	ΡŪ	$0.977 \\ 0.001$	$0.874 \\ 0.116$	$0.897 \\ 0.1$	$0.995 \\ 0.002$	$0.994 \\ 0$	0.995 0.001
Position 156 Score: 0.85	EuA u 99	EuA 00-04	$_{05-09}^{\rm EuA}$	As u 99	$\mathop{\mathrm{As}}\limits_{00-04}$	$^{\mathrm{As}}_{05-09}$	Position 158 Score: 0.87	EuA u 99	EuA 00-04	EuA 05–09	As u 99	$^{ m As}_{ m 00-04}$	$\stackrel{ m As}{_{05-09}}$
40F	$\begin{array}{c} 0.11 \\ 0.007 \\ 0.853 \end{array}$	$\begin{array}{c} 0.071 \\ 0.004 \\ 0.899 \end{array}$	$\begin{array}{c} 0.02 \\ 0.009 \\ 0.960 \end{array}$	$\begin{array}{c} 0.079 \\ 0.053 \\ 0.851 \end{array}$	$\begin{array}{c} 0.051 \\ 0.012 \\ 0.918 \end{array}$	$\begin{array}{c} 0.024 \\ 0.036 \\ 0.923 \end{array}$	AA AG GA	$\begin{array}{c} 0.936 \\ 0.009 \\ 0.02 \end{array}$	$\begin{array}{c} 0.946 \\ 0.007 \\ 0.035 \end{array}$	$\begin{array}{c} 0.963 \\ 0.012 \\ 0.011 \end{array}$	$\begin{array}{c} 0.868 \\ 0.014 \\ 0.029 \end{array}$	$\begin{array}{c} 0.943 \\ 0.019 \\ 0.01\end{array}$	$\begin{array}{c} 0.852 \\ 0.093 \\ 0.018 \end{array}$
Position 171 Score: 0.94	EuA u 99	EuA 00-04	${ m EuA}_{05-09}$	As u 99	$\mathop{\mathrm{As}}\limits_{00-04}$	$^{ m As}_{05-09}$	Position 178 Score: 0.86	EuA u 99	EuA 00-04	EuA $05-09$	As u 99	${}^{ m As}_{ m 00-04}$	$\stackrel{ m As}{_{05-09}}$
CACC	0.957 0.022	$0.98 \\ 0.012$	$0.982 \\ 0.005$	$0.913 \\ 0.041$	$0.936 \\ 0.021$	0.945 0.032	AA CC CC CC CC CC CC	$\begin{array}{c} 0.017\\ 0.620\\ 0.043\\ 0.205\\ 0.018\end{array}$	$\begin{array}{c} 0.003\\ 0.862\\ 0.073\\ 0.014\\ 0.009\end{array}$	$\begin{array}{c} 0.007\\ 0.943\\ 0.011\\ 0.007\\ 0.009\end{array}$	$\begin{array}{c} 0.124\\ 0.605\\ 0.039\\ 0.083\\ 0.047\end{array}$	$\begin{array}{c} 0.067 \\ 0.757 \\ 0.032 \\ 0.002 \\ 0.012 \end{array}$	$\begin{array}{c} 0.037 \\ 0.867 \\ 0.007 \\ 0.005 \\ 0.034 \end{array}$
Position 192 Score: 0.89	EuA u 99	EuA 00-04	$_{05-09}^{\rm EuA}$	As u 99	$\operatorname*{As}_{00-04}$	$^{\mathrm{As}}_{05-09}$	Position 201 Score: 0.99	EuA u 99	EuA 00-04	EuA 05-09	As u 99	${}^{ m As}_{ m 00-04}$	$\stackrel{ m As}{_{05-09}}$
AA- AC- AG- AT- GA-	$\begin{array}{c} 0.612\\ 0.09\\ 0.13\\ 0.025\\ 0.036\end{array}$	$\begin{array}{c} 0.538\\ 0.278\\ 0.021\\ 0.017\\ 0.032\end{array}$	$\begin{array}{c} 0.428\\ 0.38\\ 0.011\\ 0.086\\ 0.01\end{array}$	$\begin{array}{c} 0.626\\ 0.007\\ 0.089\\ 0.001\\ 0.064\end{array}$	$\begin{array}{c} 0.793 \\ 0.002 \\ 0.044 \\ 0 \\ 0 \\ 0.024 \end{array}$	$\begin{array}{c} 0.888\\ 0.006\\ 0.001\\ 0\\ 0\\ 0.009 \end{array}$	A	0.006 0.952	$0.217 \\ 0.748$	0.497 0.492	0.005 0.973	0.002 0.983	$0 \\ 0.992$

# K. Bartoszek, P. Liò, A. Sorathiya

.s As -04 05-09	76 0.044 76 0.868 44 0.085	.s As -04 05-09	03 0.052 84 0.898	.s As -04 05-09	36 0.96 25 0.001	.s As -04 05-09	19 0.961 65 0.018 12 0.013	.s As -04 05-09	02 0.106 98 0.894	s As -04 05-09	$\begin{array}{ccc} 01 & 0.004 \\ 99 & 0.995 \end{array}$	.s As -04 05-09	75 0.836 24 0.164	s As -04 05-09	6.0 66
As Ai u 99 00-	$\begin{array}{cccc} 0.093 & 0.0'\\ 0.884 & 0.8'\\ 0.023 & 0.0' \end{array}$	As A; u 99 00-	0.098 0.0 0.805 0.8	As A <sub>i</sub> u 99 00-	0.818 0.93 0.14 0.03	As Ai u 99 00-	$\begin{array}{cccc} 0.79 & 0.9 \\ 0.051 & 0.00 \\ 0.146 & 0.0 \\ 0.005 & 0 \end{array}$	As A <sub>i</sub> u 99 00-	0.023 0.00 0.974 0.99	As A <sub>i</sub> u 99 00-	0.006 0.00	As Ai u 99 00-	0.919 0.9' 0.081 0.0	As A: u 99 00-	66.0 666.0
EuA $05-09$	$\begin{array}{c} 0.018 \\ 0.963 \\ 0.008 \end{array}$	EuA $05-09$	$0.014 \\ 0.958$	EuA 05–09	0.969 0.009	EuA $05-09$	$\begin{array}{c} 0.744 \\ 0.009 \\ 0.015 \\ 0.229 \end{array}$	EuA 05-09	$0.549 \\ 0.450$	EuA $05-09$	$0.092 \\ 0.907$	EuA 05-09	$0.993 \\ 0.006$	${ m EuA}_{05-09}$	0.969
EuA 00-04	$\begin{array}{c} 0.083 \\ 0.894 \\ 0 \end{array}$	EuA 00-04	$0.012 \\ 0.960$	EuA 00-04	$0.949 \\ 0.037$	EuA 00-04	$\begin{array}{c} 0.938 \\ 0.006 \\ 0.029 \\ 0.024 \end{array}$	EuA 00-04	0.363 0.636	EuA 00-04	$0.222 \\ 0.778$	EuA 00-04	$0.994 \\ 0.006$	EuA 00-04	0 991
EuA u 99	$\begin{array}{c} 0.126 \\ 0.856 \\ 0 \end{array}$	EuA u 99	$0.042 \\ 0.903$	EuA u 99	$0.912 \\ 0.071$	EuA u 99	$\begin{array}{c} 0.686 \\ 0.014 \\ 0.285 \\ 0.012 \end{array}$	EuA u 99	0.148 0.852	EuA u 99	0.083 0.906	EuA u 99	$0.978 \\ 0.022$	EuA u 99	0 955
Position 262 Score: 0.86	- T A -	Position 319 Score: 0.87	GA TC	Position 340 Score: 0.94	CA	Position 408 Score: 0.88	AAG CAG GAG TAG	Position 457 Score: 0.96	LT CT	Position 504 Score: 0.99	υF	Position 512 Score: 0.95	AG	Position 555 Score: 0.99	0
$^{ m As}_{ m 05-09}$	$\begin{array}{c} 0.869 \\ 0.039 \\ 0.0589 \end{array}$	$^{ m As}_{05-09}$	$0 \\ 0.975$	${}^{ m As}_{ m 05-09}$	$0.916 \\ 0.084$	$^{ m As}_{05-09}$	$0.089 \\ 0.911$	$^{ m As}_{05-09}$	$\begin{array}{c} 0.121 \\ 0.879 \\ 0.894 \\ 0.106 \end{array}$	$^{ m As}_{05-09}$	0.0909	$^{ m As}_{05-09}$	$0.787 \\ 0.213$	$^{ m As}_{05-09}$	0.285
$\stackrel{ m As}{_{00-04}}$	$\begin{array}{c} 0.857 \\ 0.053 \\ 0.062 \end{array}$	${}^{ m As}_{ m 00-04}$	0.023 0.932	$\stackrel{ m As}{_{00-04}}$	$0.932 \\ 0.068$	$\stackrel{ m As}{_{00-04}}$	0.043 0.957	$\stackrel{ m As}{_{00-04}}$	$\begin{array}{c} 0.152 \\ 0.847 \\ 0.998 \\ 0.002 \end{array}$	$\stackrel{ m As}{_{00-04}}$	$0.043 \\ 0.956$	$\stackrel{ m As}{_{00-04}}$	$0.998 \\ 0.002$	$\mathop{\mathrm{As}}_{00-04}$	0 133
As u 99	$\begin{array}{c} 0.778 \\ 0.039 \\ 0.119 \end{array}$	As u 99	$0.029 \\ 0.934$	As u 99	$0.898 \\ 0.102$	As u 99	$0.04 \\ 0.96$	As u 99	$\begin{array}{c} 0.12 \\ 0.877 \\ 0.974 \\ 0.023 \end{array}$	As u 99	$0.023 \\ 0.976$	As u 99	$0.997 \\ 0.002$	As u 99	0.004
EuA 05–09	$\begin{array}{c} 0.961 \\ 0.005 \\ 0.018 \end{array}$	$E_{uA}$ 05–09	$\begin{array}{c} 0.005 \\ 0.984 \end{array}$	EuA 05–09	$0.989 \\ 0.005$	EuA 05–09	$0.01 \\ 0.99$	EuA 05–09	$\begin{array}{c} 0.021 \\ 0.979 \\ 0.450 \\ 0.549 \end{array}$	$E_{uA}$ 05–09	$\begin{array}{c} 0.01 \\ 0.989 \end{array}$	EuA 05–09	0.96 0.039	EuA 05–09	0.155
EuA 00-04	$\begin{array}{c} 0.933 \\ 0.009 \\ 0.045 \end{array}$	EuA 00-04	$0.067 \\ 0.924$	EuA 00-04	$0.984 \\ 0.011$	EuA 00-04	0.005 0.991	EuA 00-04	$\begin{array}{c} 0.041 \\ 0.947 \\ 0.636 \\ 0.363 \end{array}$	EuA 00-04	$\begin{array}{c} 0.006 \\ 0.994 \end{array}$	EuA 00-04	$0.999 \\ 0.001$	EuA 00-04	0.06
EuA u 99	$\begin{array}{c} 0.871 \\ 0.016 \\ 0.075 \end{array}$	EuA u 99	$0.09 \\ 0.887$	EuA u 99	$0.967 \\ 0.015$	EuA u 99	0.004 0.996	EuA u 99	$\begin{array}{c} 0.045\\ 0.946\\ 0.852\\ 0.148\end{array}$	EuA u 99	$\begin{array}{c} 0.002 \\ 0.998 \end{array}$	EuA u 99	$0.998 \\ 0$	EuA u 99	0.009
Position 202 Score: 0.88	-1 U (	Position 271 Score: 0.96	— —	Position 328 Score: 0.95	υH	Position 403 Score: 0.96	AT GT	Position 433 Score: 0.85	LT C C A	Position 497 Score: 0.98	A G	Position 507 Score: 0.996	A Q	Position 548 Score: 0.93	AC.

# Influenza Differentiation and Evolution

${}^{ m As}_{05-09}$	0.035 0.945	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.055 \\ 0.805 \\ 0.14 \end{array}$	$\mathbf{As}_{05-09}$	0.898 0.096	$^{\mathrm{As}}_{05-09}$	$0.993 \\ 0.007$	$^{\mathrm{As}}_{05-09}$	$0.922 \\ 0.043$	$^{ m As}_{05-09}$	$\begin{array}{c} 0.023 \\ 0.941 \\ 0.035 \end{array}$	$^{ m As}_{05-09}$	$0.956 \\ 0.014$
${ m As}_{00-04}$	0.02 0.936	${}^{ m As}_{ m 00-04}$	$\begin{array}{c} 0.064 \\ 0.935 \\ 0 \end{array}$	${}^{ m As}_{ m 00-04}$	$0.983 \\ 0.013$	${}^{ m As}_{ m 00-04}$	$0.994 \\ 0.006$	$\stackrel{ m As}{_{00-04}}$	0.899 0.09	$\mathop{\mathrm{As}}_{00-04}$	$\begin{array}{c} 0.049 \\ 0.924 \\ 0.025 \end{array}$	$\mathop{\mathrm{As}}_{00-04}$	$0.939 \\ 0.043$
As u 99	$0.045 \\ 0.913$	As u 99	$\begin{array}{c} 0.087 \\ 0.912 \\ 0.001 \end{array}$	As u 99	$0.971 \\ 0.021$	As u 99	0.906 0.093	As u 99	$0.881 \\ 0.069$	As u 99	$\begin{array}{c} 0.039 \\ 0.88 \\ 0.081 \end{array}$	As u 99	$0.759 \\ 0.2$
${ m EuA}_{05-09}$	0.009 0.982	EuA 05-09	$\begin{array}{c} 0.013 \\ 0.906 \\ 0.077 \end{array}$	EuA 05-09	$0.988 \\ 0.01$	EuA 05-09	$0.877 \\ 0.122$	EuA $05-09$	$0.977 \\ 0.02$	EuA 05-09	0.007 0.976 0.016	EuA $05-09$	$0.574 \\ 0.418$
EuA 00-04	0.003 0.98	EuA 00-04	$\begin{array}{c} 0.017 \\ 0.978 \\ 0.002 \end{array}$	EuA 00-04	$\begin{array}{c} 0.977 \\ 0.021 \end{array}$	EuA 00-04	$0.822 \\ 0.176$	EuA 00-04	$0.976 \\ 0.018$	EuA 00-04	$\begin{array}{c} 0.035 \\ 0.948 \\ 0.015 \end{array}$	$_{00-04}^{\rm EuA}$	$0.754 \\ 0.239$
EuA u 99	0.032 0.955	EuA u 99	$\begin{array}{c} 0.042 \\ 0.955 \\ 0.001 \end{array}$	EuA u 99	$0.98 \\ 0.017$	EuA u 99	$0.923 \\ 0.076$	EuA u 99	0.983 0.015	EuA u 99	$\begin{array}{c} 0.021 \\ 0.913 \\ 0.043 \end{array}$	EuA u 99	0.096
Position 626 Score: 0.95	AGA TGG	Position 698 Score: 0.91	AC GA GG	Position 704 Score: 0.92	AT GT	Position 788 Score: 0.95	GATGG GACGG	Position 830 Score: 0.89	AT GT	Position 865 Score: 0.89	AG A- G-	Position 882 Score: 0.85	ΥÜ
$^{ m As}_{05-09}$	$\begin{array}{c} 0.052 \\ 0.035 \\ 0.027 \\ 0.863 \end{array}$	$^{ m As}_{05-09}$	$\begin{array}{c} 0.893 \\ 0.052 \\ 0.055 \end{array}$	$^{ m As}_{05-09}$	$\begin{array}{c} 0.878 \\ 0.099 \\ 0.0237 \end{array}$	$^{ m As}_{05-09}$	$\begin{array}{c} 0.932 \\ 0.035 \\ 0.0292 \end{array}$	$^{ m As}_{05-09}$	$\begin{array}{c} 0.966 \\ 0.01 \\ 0.022 \end{array}$	$^{ m As}_{05-09}$	0.95 0.025	$^{ m As}_{05-09}$	$\begin{array}{c} 0\\ 0.949\\ 0.03\\ 0.012 \end{array}$
${}^{ m As}_{ m 00-04}$	$\begin{array}{c} 0.03 \\ 0.021 \\ 0.053 \\ 0.833 \end{array}$	${}^{ m As}_{ m 00-04}$	$\begin{array}{c} 0.917 \\ 0.019 \\ 0.064 \end{array}$	${}^{ m As}_{ m 00-04}$	$\begin{array}{c} 0.898 \\ 0.051 \\ 0.051 \end{array}$	${}^{ m As}_{ m 00-04}$	$\begin{array}{c} 0.907 \\ 0.011 \\ 0.04 \end{array}$	$\stackrel{ m As}{_{00-04}}$	$\begin{array}{c} 0.929 \\ 0.019 \\ 0.051 \end{array}$	$\stackrel{ m As}{_{00-04}}$	$0.888 \\ 0.071$	$\mathop{\mathrm{As}}_{00-04}$	$\begin{array}{c} 0\\ 0.898\\ 0.059\\ 0.011\end{array}$
As u 99	$\begin{array}{c} 0.098 \\ 0.043 \\ 0.041 \\ 0.772 \end{array}$	As u 99	$\begin{array}{c} 0.832 \\ 0.071 \\ 0.087 \end{array}$	As u 99	$\begin{array}{c} 0.826 \\ 0.135 \\ 0.039 \end{array}$	As u 99	$\begin{array}{c} 0.897 \\ 0.047 \\ 0.013 \end{array}$	As u 99	$\begin{array}{c} 0.914 \\ 0.047 \\ 0.039 \end{array}$	As u 99	$0.857 \\ 0.105$	As u 99	$\begin{array}{c} 0\\ 0.907\\ 0.053\\ 0.025\end{array}$
EuA 05–09	$\begin{array}{c} 0.014 \\ 0.007 \\ 0.008 \\ 0.953 \end{array}$	EuA 05–09	$\begin{array}{c} 0.968 \\ 0.012 \\ 0.018 \end{array}$	EuA 05–09	$\begin{array}{c} 0.972 \\ 0.02 \\ 0.008 \end{array}$	EuA 05–09	$\begin{array}{c} 0.977\\ 0.009\\ 0.011\end{array}$	EuA 05–09	$\begin{array}{c} 0.482 \\ 0.508 \\ 0.008 \end{array}$	EuA 05–09	$0.952 \\ 0.021$	EuA 05–09	$\begin{array}{c} 0.0968 \\ 0.857 \\ 0.011 \\ 0.035 \end{array}$
EuA 00-04	$\begin{array}{c} 0.012 \\ 0.012 \\ 0.035 \\ 0.921 \end{array}$	EuA 00-04	$\begin{array}{c} 0.966 \\ 0.014 \\ 0.02 \end{array}$	EuA 00-04	0.939 0.026 0.035	EuA 00-04	$\begin{array}{c} 0.951 \\ 0.01 \\ 0.036 \end{array}$	EuA 00-04	$\begin{array}{c} 0.698 \\ 0.267 \\ 0.035 \end{array}$	EuA 00-04	600.0 676.0	EuA 00-04	$\begin{array}{c} 0.103 \\ 0.845 \\ 0.026 \\ 0.024 \end{array}$
EuA u 99	$\begin{array}{c} 0.042 \\ 0.022 \\ 0.044 \\ 0.834 \end{array}$	EuA u 99	$\begin{array}{c} 0.925 \\ 0.029 \\ 0.044 \end{array}$	EuA u 99	$\begin{array}{c} 0.898 \\ 0.058 \\ 0.044 \end{array}$	EuA u 99	$\begin{array}{c} 0.929 \\ 0.021 \\ 0.038 \end{array}$	EuA u 99	$\begin{array}{c} 0.893 \\ 0.063 \\ 0.044 \end{array}$	EuA u 99	$0.891 \\ 0.094$	EuA u 99	$\begin{array}{c} 0\\ 0.944\\ 0.048\\ 0.009\end{array}$
Position 560 Score: 0.87	AA AG GT TC	Position 677 Score: 0.86	ΤCΑ	Position 701 Score: 0.87	ΑΩΗ	Position 761 Score: 0.99	AA CA GA	Position 819 Score: 0.85	AC AT CA	Position 854 Score: 0.9	-AG	Position 876 Score: 0.85	40H

K. Bartoszek, P. Liò, A. Sorathiya

444

$\begin{array}{ccc} & As & As \\ 9 & 00-04 & 05-09 \end{array}$	81 0.023 0.036 53 0.053 0.044 81 0.913 0.898	9 00-04 05-09	37 0.082 0.072 49 0.012 0.035 76 0.855 0.869 39 0.051 0.023	As As As 9 00-04 05-09	06 0.9 0.886 38 0.068 0.056	As As As 9 00-04 05-09	07 0.004 0.084 77 0.963 0.907	As As As 9 00-04 05-09	24 0.863 0.904 12 0.02 0.012 34 0.098 0.067	As As As 9 00-04 05-09	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	, As As 9 00-04 05-09	$\begin{array}{cccccccccccccccccccccccccccccccccccc$
uA As -09 u 9	$\begin{array}{cccc} 011 & 0.08 \\ 014 & 0.06 \\ 972 & 0.83 \end{array}$	uA As -09 u 9	022 0.13 009 0.04 961 0.77 008 0.03	uA As -09 u 9	967 0.80 018 0.08	uA As -09 u 9	051 0.00 947 0.97	uA As -09 u 9	969 0.82 019 0.04 007 0.08	uA As -09 u 9	0 0.00 1186 0.00 005 0.08 3388 0 335 0 0 335 0 0 0 0 0 0 0 0 0 0 0 0 0 0	uA As -09 u 9	011 0.1- 988 0.8
3uA E 0-04 05	.009 0.1 .017 0.1 .966 0.1	EuA E 0-04 05	.023 0.1 0.01 0.1 .933 0.1 .0347 0.1	3uA E 0-04 05	.966 0.1	3uA E 0-04 05	.015 0. .983 0.	3uA E 0-04 05	0.94 0.1 0.02 0.1 0.035 0.1	3uA E 0-04 05	0 0.15 0 0.07 0.0 0.82 0.1 0.25 0.1 104 0.1 0.015 0.31 0.1	3uA E 0-04 05	.007 0. .992 0.
EuA ] u 99 0	0.039 C 0.034 C 0.921 C	EuA 1 u 99 0	0.065 C 0.022 0.869 C 0.044 0	EuA ] u 99 0	0.912 C 0.044 C	EuA ] u 99 0	0.015 C 0.975 C	EuA ] u 99 0	0.89 0.048 0.041 C	EuA ] u 99 0	$\begin{array}{c} 0\\ 0.061\\ 0.041\\ 0.059\\ 0.011\\ 0.015\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\$	EuA ] u 99 0	0.032 C 0.96 C
Position 975 Score: 0.88	AC GC TC	Position 1053 Score: 0.85	AG AG GG TC	Position 1086 Score: 0.86	A G	Position 1161 Score: 0.92	GG	Position 1173 Score: 0.87	TGGAT TGGCT TGGTT	Position 1196 Score: 0.86	GGATTCGCGGGTCAGGCTATGAGACA GAATCACGCTCAGGATATGAGACA GAATCACGCTCAGGTTATGAGAACC GAATTACGCTCAGGTTATGAGAACT TAATTCCAGGAGCGGCTTTGAAACG TAATTCCAGGAGCGGCTTTGAAACG GGATTCCAGGAGCGGCTTTGAGAACG GAATTCAAGAAACGGTTTGAGAACG GAAGTCAAGGAACGGTTATGAGAACG GAAGTCAACGACCTTAGGGTATGAGAACG CAGACTTAGAAAGGGGTTTGAGAACG	Position 1249 Score: 0.94	A C
$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.914 \\ 0.062 \\ 0.024 \end{array}$	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.923 \\ 0.041 \\ 0.04 \end{array}$	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.024 \\ 0.903 \\ 0.072 \end{array}$	$^{\mathrm{As}}_{05-09}$	$0.005 \\ 0.995$	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.816 \\ 0.104 \end{array}$	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.149\\ 0\\ 0.004\\ 0.118\\ 0.215\\ 0.082\\ 0.1\\ 0.1\\ 0.064\end{array}$	$_{05-09}^{\rm As}$	$0.035 \\ 0.945$
$_{00-04}^{\rm As}$	$\begin{array}{c} 0.887 \\ 0.062 \\ 0.051 \end{array}$	${}^{ m As}_{00-04}$	$\begin{array}{c} 0.886 \\ 0.09 \\ 0.023 \end{array}$	${\rm As}_{00-04}$	$\begin{array}{c} 0.042 \\ 0.886 \\ 0.062 \end{array}$	${\rm As}_{00-04}$	$0.001 \\ 0.999$	${\rm As}_{00-04}$	$\begin{array}{c} 0.818 \\ 0.076 \end{array}$	${\rm As}_{00-04}$	$\begin{array}{c} 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0.036\\ 0.464 \end{array}$	$\mathop{\rm As}_{00-04}$	$0.014 \\ 0.936$
$^{\rm As}_{ m u}$ 99	$\begin{array}{c} 0.793 \\ 0.168 \\ 0.039 \end{array}$	$^{\mathrm{As}}_{\mathrm{u}\ 99}$	$\begin{array}{c} 0.747 \\ 0.164 \\ 0.081 \end{array}$	As u 99	$\begin{array}{c} 0.049 \\ 0.81 \\ 0.141 \end{array}$	As u 99	0.006 0.993	As u 99	$0.709 \\ 0.131$	As u 99	$egin{array}{c} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 $	$^{ m As}_{ m u}$ 99	$0.038 \\ 0.912$
EuA 05 $-09$	$\begin{array}{c} 0.974 \\ 0.018 \\ 0.008 \end{array}$	EuA 05-09	$\begin{array}{c} 0.963 \\ 0.018 \\ 0.011 \end{array}$	EuA 05-09	$\begin{array}{c} 0.011 \\ 0.97 \\ 0.018 \end{array}$	EuA 05-09	$\begin{array}{c} 0.088 \\ 0.911 \end{array}$	EuA 05-09	$0.941 \\ 0.017$	EuA 05-09	$\begin{array}{c} 0.001\\ 0.303\\ 0.02\\ 0.07\\ 0.032\\ 0.024\\ 0.105\end{array}$	EuA 05-09	$0.006 \\ 0.981$
EuA 00-04	$\begin{array}{c} 0.952 \\ 0.013 \\ 0.035 \end{array}$	$E_{uA}$ 00-04	$\begin{array}{c} 0.957 \\ 0.026 \\ 0.009 \end{array}$	EuA 00-04	$\begin{array}{c} 0.046 \\ 0.938 \\ 0.016 \end{array}$	EuA 00-04	$\begin{array}{c} 0.049 \\ 0.951 \end{array}$	EuA 00-04	$0.95 \\ 0.01$	EuA 00-04	$\begin{array}{c} 0 \\ 0.201 \\ 0.0307 \\ 0.001 \\ 0.005 \\ 0 \\ 0 \\ 0 \\ 0.223 \end{array}$	EuA 00-04	$0.012 \\ 0.98$
EuA u 99	$\begin{array}{c} 0.908 \\ 0.048 \\ 0.044 \end{array}$	EuA u 99	$\begin{array}{c} 0.893 \\ 0.056 \\ 0.039 \end{array}$	EuA u 99	$\begin{array}{c} 0.061 \\ 0.882 \\ 0.056 \end{array}$	EuA u 99	$0.009 \\ 0.991$	EuA u 99	$\begin{array}{c} 0.731 \\ 0.167 \end{array}$	EuA u 99	$egin{array}{c} 0 \\ 0.01 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0.12 \end{array}$	EuA u 99	0.022 0.955
Position 942 Score: 0.88	ΥÜΗ	Position 1002 Score: 0.85	AA GA TC	Position 1084 Score: 0.86	Ϋ́́ύΕ	Position 1137 Score: 0.99	AT GT	Position 1164 Score: 0.86	AATG GATG	Position 1182 Score: 0.86	AGAACTAAAAGTAACAGACTTAGAA AGAACCAAAAGCACTAATTCCAGGA AGGACCAAAAGCACTAAGGACTTAGAA AGGACCAAAAGTAACAGACTTAGAA AGGACCAAAAGTAACAGAGACTTAGAA AGAACAATCAGCGAGAAGTCACGCT AGAACAATCAGCAGGGAAGTTCGCGGCT AGAACAATCAGCAGGAGAAGTTCGCGGT AGAACTAAAAGCATTAGTTCGAGGAT AGAACTAAAAGCATTAGTTCAAGAGA	Position 1239 Score: 0.93	GAA

$^{\mathrm{As}}_{05-09}$	$0 \\ 0.982$	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.039\\ 0.006\\ 0.019\\ 0.937\end{array}$	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.915 \\ 0.035 \\ 0.029 \end{array}$	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.012 \\ 0.988 \end{array}$	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.013 \\ 0.987 \end{array}$		
${}^{ m As}_{ m 00-04}$	0 0.988	${}^{ m As}_{ m 00-04}$	$\begin{array}{c} 0.015 \\ 0.033 \\ 0.017 \\ 0.935 \end{array}$	$\stackrel{ m As}{_{00-04}}$	$\begin{array}{c} 0.885 \\ 0.021 \\ 0.051 \end{array}$	${}^{ m As}_{ m 00-04}$	$\begin{array}{c} 0.001 \\ 0.998 \end{array}$	$\stackrel{ m As}{_{00-04}}$	$\begin{array}{c} 0.002 \\ 0.998 \end{array}$		
As u 99	$0 \\ 0.959$	As u 99	$\begin{array}{c} 0.031 \\ 0.031 \\ 0.075 \\ 0.864 \end{array}$	As u 99	$\begin{array}{c} 0.838 \\ 0.058 \\ 0.075 \end{array}$	As u 99	$0.001 \\ 0.997$	As u 99	$\begin{array}{c} 0.004 \\ 0.994 \end{array}$		
EuA 05–09	0.036 0.946	EuA 05–09	$\begin{array}{c} 0.011 \\ 0.024 \\ 0.002 \\ 0.963 \end{array}$	EuA 05–09	$\begin{array}{c} 0.966 \\ 0.012 \\ 0.016 \end{array}$	EuA 05–09	$\begin{array}{c} 0.061 \\ 0.938 \end{array}$	EuA 05-09	$\begin{array}{c} 0.068 \\ 0.932 \end{array}$		
EuA 00-04	$\begin{array}{c} 0.117 \\ 0.845 \end{array}$	EuA 00-04	$\begin{array}{c} 0.009\\ 0.047\\ 0.002\\ 0.942\end{array}$	EuA 00-04	$\begin{array}{c} 0.925 \\ 0.021 \\ 0.036 \end{array}$	EuA 00-04	$0.134 \\ 0.866$	EuA 00-04	$0.143 \\ 0.855$		
EuA u 99	$0.016 \\ 0.923$	EuA u 99	$\begin{array}{c} 0.021 \\ 0.045 \\ 0.021 \\ 0.913 \end{array}$	EuA u 99	$\begin{array}{c} 0.896 \\ 0.034 \\ 0.052 \end{array}$	EuA u 99	$0.02 \\ 0.979$	EuA u 99	$0.027 \\ 0.961$		
Position 1300 Score: 0.98	CGTC CGTC	Position 1331 Score: 0.86	400H	Position 1464 Score: 0.86	GCA	Position 1490 Score: 0.99	GA	Position 1497 Score: 0.98	0		
$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.913 \\ 0.023 \\ 0.064 \end{array}$	$^{\mathrm{As}}_{05-09}$	$\begin{array}{c} 0.899\\ 0.053\\ 0.047\end{array}$	$\overset{\mathrm{As}}{_{05-09}}$	0.86 0.06 0.05	$\overset{\mathrm{As}}{_{05-09}}$	0.036	$^{\mathrm{As}}_{05-09}$	$0.013 \\ 0.987$	$A_{\rm S}^{\rm AS}$ 05–09	$\begin{array}{c} 0.007\\ 0.882\\ 0.001\\ 0.006\\ 0.006\\ 0.002\\ 0.016\end{array}$
${}^{ m As}_{ m 00-04}$	$\begin{array}{c} 0.909 \\ 0.051 \\ 0.04 \end{array}$	${}^{ m As}_{ m 00-04}$	0.938 0.029 0.033	$\stackrel{ m As}{_{00-04}}$	$\begin{array}{c} 0.882 \\ 0.074 \\ 0.029 \end{array}$	${}^{ m As}_{00-04}$	$0.021 \\ 0.93$	$\stackrel{ m As}{_{00-04}}$	$0.001 \\ 0.990$	$\stackrel{ m As}{_{00-04}}$	$\begin{array}{c} 0.115\\ 0.555\\ 0.001188\\ 0.129\\ 0.005\\ 0.0167\end{array}$
As u 99	$\begin{array}{c} 0.839 \\ 0.037 \\ 0.124 \end{array}$	As u 99	$\begin{array}{c} 0.834 \\ 0.06 \\ 0.105 \end{array}$	As u 99	$\begin{array}{c} 0.814 \\ 0.119 \\ 0.047 \end{array}$	As u 99	$\begin{array}{c} 0.037 \\ 0.91 \end{array}$	As u 99	$0.003 \\ 0.995$	As u 99	$\begin{array}{c} 0.058\\ 0.668\\ 0\\ 0.035\\ 0.012\\ 0.042 \end{array}$
EuA 05-09	$\begin{array}{c} 0.975 \\ 0.007 \\ 0.018 \end{array}$	EuA 05-09	$\begin{array}{c} 0.972 \\ 0.009 \\ 0.019 \end{array}$	EuA 05-09	$\begin{array}{c} 0.952 \\ 0.012 \\ 0.006 \end{array}$	EuA 05-09	$\begin{array}{c} 0.01 \\ 0.973 \end{array}$	EuA 05-09	$\begin{array}{c} 0.062 \\ 0.921 \end{array}$	EuA 05-09	$\begin{array}{c} 0.004\\ 0.841\\ 0.018\\ 0.01\\ 0.042\\ 0.016\\ 0.016\end{array}$
EuA 00-04	$\begin{array}{c} 0.95 \\ 0.035 \\ 0.015 \end{array}$	EuA 00-04	$\begin{array}{c} 0.968 \\ 0.017 \\ 0.015 \end{array}$	EuA 00-04	$\begin{array}{c} 0.918 \\ 0.039 \\ 0.01 \end{array}$	EuA 00-04	$\begin{array}{c} 0.015 \\ 0.958 \end{array}$	EuA 00-04	$0.142 \\ 0.854$	EuA 00-04	$\begin{array}{c} 0.025\\ 0.743\\ 0.074\\ 0.016\\ 0.039\\ 0.021\end{array}$
EuA u 99	$\begin{array}{c} 0.909\\ 0.036\\ 0.055\end{array}$	EuA u 99	$\begin{array}{c} 0.915 \\ 0.032 \\ 0.05 \end{array}$	EuA u 99	$\begin{array}{c} 0.865 \\ 0.065 \\ 0.029 \end{array}$	EuA u 99	$0.017 \\ 0.957$	EuA u 99	$0.027 \\ 0.951$	EuA u 99	$\begin{array}{c} 0.036\\ 0.636\\ 0.074\\ 0.047\\ 0.031\\ 0.018 \end{array}$
Position 1252 Score: 0.88	400	Position 1317 Score: 0.87	ΥÜΗ	Position 1460 Score: 0.88	AG CA GG	Position 1469 Score: 0.9	GG TG	Position 1493 Score: 0.97	66	Position 1546 Score: 0.86	AAAGTGCTTGTTTCTACT AAACTGCTTGTTTCTAC AAACTGCTTGTTTCTACT GTTTCTACT TGTTTCTACT

### 4.5. Column entropy

We also looked at how variable all of our gene segments were. To do this we calculated the entropy of each site in the alignment. Sites with gaps were excluded. The formula used to estimate the entropy for the j-th column is naturally,

$$\hat{I}_j = -\sum_{R \in \{A,C,G,T\}} \frac{n_R}{n} \log \frac{n_R}{n} \,, \tag{1}$$

where n is the amount of sequences, and  $n_R$  is the count of the given residue in the considered column. This famous formula can be thought to represent how "random" the column is. If the column is constant it will be 0 and if all residues are equally present (*i.e.* for all  $R n_R/n = 0.25$ ) then it will obtain its maximum of log 4. The results are in Figs. 6, 7 and 8.



Fig. 6. Estimated entropies at individual sites in alignment of HA gene segment sequences and compared with breaking the sequences according to hosts and geographical regions.

What can be seen from the figure is that there is not that much difference between hosts and regions in each gene segment but there are differences between the gene segments. This could be because of different selective pressures on the different gene segments due to the M protein being an internal protein while the HA and NA are external proteins. External proteins in viruses are more liable for change as they are responsible for the virus having the ability to infect. Therefore the difference between gene segments is rather obvious but the lack of difference between hosts and regions could be interesting for further biological and modelling insight.



Fig. 7. Estimated entropies at individual sites in alignment of M gene segment sequences and compared with breaking the sequences according to hosts and geographical regions. Entropy not calculated for columns with gaps.



Fig. 8. Estimated entropies at individual sites in alignment of NA gene segment sequences and compared with breaking the sequences according to hosts and geographical regions. Entropy not calculated for columns with gaps.

### 5. Discussion

Our aim was to using a simple methodology to find regions of up to 25 base pairs in influenza viral gene segment sequences that differentiate between certain factors (host, geography, time period, Hx strain). We wanted to see whether anything novel could be found using a rather "rough look" at data. The aim was achieved in a sense that a number of promising places in the gene segments were found. Because the methodology is very simple and requires some heuristic scoring mechanism its drawbacks have to be discussed.

### 5.1. Sequence acquisition

The sequences were downloaded from GISAID [4]. Therefore the methodology is dependent on us obtaining a random sample of sequences which when looking at the sequence breakdown can be seen to be immediately violated. It is clear that the majority of sequences have been collected from the past ten years and from human hosts. Another issue is that sometimes the lengths of the downloaded sequence were drastically different. For example in the HA gene segment of H1N1 about 1000 sequences were less than half the length of the others (so one third of the sequences). They had to be discarded as they made the alignment unsure and also the program was finding only differences between these two groups *i.e.* gaps/no gaps. It cannot be assessed therefore whether this could have caused any bias. One issue with the downloaded sequences is that we do not have all the information for every single one, sometimes region, host or time is missing. This then cuts down the sample sizes in some parts of the study.

#### 5.2. Method

The first step of the analysis is to align the acquired sequences. This was done in Clustal 10.2 [5]. Due to the huge number of sequences (a couple of thousand in each case) the alignment has to be done approximately. Visual inspection of it showed extremely high similarity between all of the sequences and no clear misaligned regions could be seen. But this is the grand picture, small regions/individual bases could be very well misaligned which can have profound effects on the method of searching for characteristic differences. It was noticed very often that in most places one had a huge number (around 200) of different values of the metasubsequence that had one, two, less than ten sequences. This could be due to mutations in those individual sequences (which could very well be important) but just as well due to a misalignment. Another observation made is that especially in the M and NA gene segments nearly every single metasubsequence of length 25 was seen by the contingency table test (at a *p*-value of  $1^{-183}$ ) to be significant in the spatial and temporal analysis. Such a huge number of positions has to be cut down naturally. At the moment this is done first by only considering those values of the metasubsequence that have at least 2% of the sequences altogether and then assigning scores to each contingency table and we are left with the presented results. Because of the non-independence the *p*-values should be treated as a cut-off value and not understood in the statistical sense. A last issue is how to combine factors, *e.g.* countries, time — periods. The division of time-periods should be sequential. Here we divided the time into three periods to keep the number of sequences in each bin approximately the same. But this might cause us to miss some effects inside these bins and there also comes the potential problem should the bins be the same for every geographical region. An issue in the creation of categories is computational. for this to be feasible there cannot be too many, we considered a maximum of 12 categories (with approximately 10000 sequences). Our analysis does not address these issues. It is not its aim. The aim is to find as many as possible potentially interesting sites in the flu gene segments that show differentiation between hosts, geography and time by a simple first glance.

### 5.2.1. Scoring method

To cut down on the huge number of significant tables we devised a scoring method described in Section 3.3. The intuition behind it is that it will highly score those tables where given the group we get a lot of information about the sequence. An obvious issue with the method is to decide the cut-off value. At the moment we choose a level such that the number of results would be manageable to look through. When looking at Table XIII we can see that positions 201 and 202 are in it separately. This is because together they did not generate a significant contingency table. Something else that can be noticed in the presented results is that they show positions which start and end with constant columns. These of course do not influence the score but they might play a role in the *p*-value. We have to remember that to calculate the score we remove those rows that have less than 2% of the sequences. There could be changes on these "constant" positions in the removed rows which decide that the table is significant.

# 5.3. Evolution model

In a recently published paper [16] the authors present a mathematical model which they claim well describes the creation of viral phylogenetic trees. From our perspective when analyzing the influenza differentiating metasubsequences, the important conclusion from [16] is that one small change can create a virus version which will be present in significant numbers. This is strikingly visible (unless this is due to a sampling bias over which we have no control) when looking at the values of the metasubsequence at position 1023 of the M gene segment. We have a version of the virus that in the Americas was very scarce and then recently started to become visible in a huge proportion. When we looked at the contingency tables in every single sequence of positions in the alignment we observed that there could be around 200 realizations in those positions that had 1, 2, 3 below 10 sequences having that value of the metasubsequence. We said that this could be due to a misalignment but in view of the model and the result in the M gene segment they could just as well be real changes that might or might not become dominant. At the moment this was noticed only in one place.

### 6. Further developments

The analysis done here is a very superficial first glance at differentiation between influenza viruses due to factors like host, geography or time period. Interactions between time and region were also looked at. The next steps should be introducing models of strain evolution, especially to consider the phylogenetic dependence between them.

K.B. is partially funded by the Kungliga Vetenskapsakademien.

#### REFERENCES

- N. Voirin, B. Barret, M.-H. Metzger, P. Vanhems, Journal of Hospital Infection 71, 1 (2009).
- [2] T.P. Weber, N.I. Stilianakis, Journal of Infection 57, 361 (2008).
- [3] N. Ferguson, *Nature* **446**, 733 (2007).
- [4] GISAID Platform, http://platform.gisaid.org/
- [5] M. Larkin et al., Bioinformatics 23, 2947 (2007).
- [6] W. Ewens, G. Grant, Statistical Methods in Bioinformatics, Springer, New York 2005.
- [7] A. Tamuri, M. dos Reis, A. Hay, R. Goldstein, Identifying Changes in Selective Constraints: Host Shifts in Influenza, PLoS Computational Biology 5(11) (2009).
- [8] Y. Furuse, A. Suzuki, T. Kamigaki, H. Oshitani, Evolution of the M Gene of the Influenza A Virus in Different Host Species: Large Scale Sequence Analysis, Virology Journal 6 (2009).
- [9] O. Miotto, A. Heiny, T. Tan, J. August, V. Brusic, Identification of Humanto-human Transmissibility Factors in PB2 Proteins of Influenza a by Largescale Mutual Information Analysis, BMC Bioinformatics 9(Suppl 1) (2008).

- [10] W. Pirovano, K. Feenstra, J. Herings, Nucleic Acids Research 34, 6540 (2006).
- [11] K. Feenstra, W. Pirovano, K. Krab, J. Herings, Nucleic Acids Research 35, W495 (2007).
- [12] Perl v5.8.5., http://www.perl.com
- [13] J. Stajich et al., Genome Research 12, 1161 (2002).
- [14] M. Kospach, Statistics-Distributions-1.02, http://search.cpan.org/ mikek/Statistics-Distributions-1.02/
- [15] D. Huson, D. Richter, C. Rausch, T.Dezulian, M. Franz, R Rupp, BMC Bioinformatics 8, 460 (2007).
- [16] T. Liggett, R. Schinazi, Journal of Applied Probability 46, 601 (2009).